

Two Heads Are Better Than One: Improving Fake News Video Detection by Correlating with Neighbors

Peng Qi^{1,2}, Yuyang Zhao³, Yufeng Shen², Wei Ji³, Juan Cao^{1,2*} and Tat-Seng Chua³

¹ Key Laboratory of Intelligent Information Processing,
Institute of Computing Technology, Chinese Academy of Sciences

² University of Chinese Academy of Sciences

³ National University of Singapore

{qipeng, caojuan}@ict.ac.cn, yuyang.zhao@u.nus.edu,
shenyufeng22@mails.ucas.ac.cn, {jiwei, dcscts}@nus.edu.sg

Abstract

The prevalence of short video platforms has spawned a lot of fake news videos, which have stronger propagation ability than textual fake news. Thus, automatically detecting fake news videos has been an important countermeasure in practice. Previous works commonly verify each news video individually with multimodal information. Nevertheless, news videos from different perspectives regarding the same event are commonly posted together, which contain complementary or contradictory information and thus can be used to evaluate each other mutually. To this end, we introduce a new and practical paradigm, *i.e.*, cross-sample fake news video detection, and propose a novel framework, Neighbor-Enhanced fake News Video Detection (NEED), which integrates the neighborhood relationship of new videos belonging to the same event. NEED can be readily combined with existing single-sample detectors and further enhance their performances with the proposed *graph aggregation (GA)* and *debunking rectification (DR)* modules. Specifically, given the feature representations obtained from single-sample detectors, GA aggregates the neighborhood information with the dynamic graph to enrich the features of independent samples. After that, DR explicitly leverages the relationship between debunking videos and fake news videos to refute the candidate videos via textual and visual consistency. Extensive experiments on the public benchmark demonstrate that NEED greatly improves the performance of both single-modal (up to 8.34% in accuracy) and multimodal (up to 4.97% in accuracy) base detectors. Codes are available in <https://github.com/ICTMCG/NEED>.

1 Introduction

“Listen to both sides and you will be enlightened; heed only one side and you will be benighted.”

— Zheng Wei (Tang Dynasty)

*Corresponding author.

Event: Man cries in pain as his clothing store is flooded.

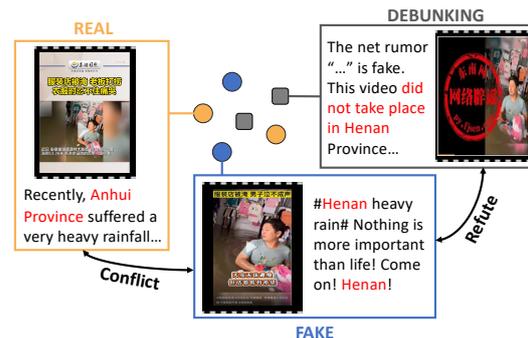


Figure 1: A set of videos belonging to the same event. Fake news videos contain conflicting information with the real ones, and the debunking videos can refute the mismatched information in the fake news videos.

The dissemination of fake news has become an important social issue which poses real-world threats to politics (Fisher et al., 2016), finance (El-Boghdady, 2013), public health (Naeem and Bhatti, 2020), *etc.* Recently, the prevalence of short video platforms has spawned a lot of fake news videos, which are more convincing and easier to spread compared to textual fake news (Sundar et al., 2021). The Cyberspace Administration of China reported that five of the seven core rumors circulating in the china eastern airlines crash incident originated from short video platforms (Cyberspace Administration of China, 2022). Statistics from another study also reveal the powerful propagation of fake news videos, which reports that only 124 TikTok fake news about COVID-19 gained more than 20 million views and 2 million likes, comments and shares, causing negative influences on millions of people (Brandy Zadrozny, 2021). Therefore, developing automatic detection techniques for fake news videos is urgent to mitigate their negative impact.

In view of the practicality of fake news video detection, previous works (Hou et al., 2019; Med-

¹[tiktok.com](https://www.tiktok.com). A popular short-form video sharing platform.

ina Serrano et al., 2020; Choi and Ko, 2021; Shang et al., 2021; Qi et al., 2023) leverage the heterogeneous multimodal information of an individual news video for corroboration. However, fake news is intentionally created to mislead consumers (Shu et al., 2017) and thus the multimodal components show few abnormalities after deliberate fabrication. In addition, fake news videos typically contain real news videos where only some frames or the textual description has been maliciously modified to alter their meanings (Qi et al., 2023). The above characteristics demonstrate that the deliberate fabrication and malicious modification are inconspicuous in a single video, leading to low effectiveness of independent detection by existing works.

In real-world scenarios, when a news event emerges, multiple related videos from different perspectives are posted, including fake news, real news, and debunking videos. Such news videos contain *complementary* or *contradictory* information, which can be used to evaluate each other mutually. As shown in Figure 1, on the one hand, the fake news video contains conflict information with the real one (*i.e.*, different locations: “Anhui” province *v.s.* “Henan” province). Furthermore, debunking videos also exist in some events, and can easily detect the corresponding fake news by providing fact-based *refutations*. In a newly released dataset (Qi et al., 2023) based on short video platforms, 54% of events containing fake news videos also have debunking videos, but 39% of events with debunking videos still had fake news videos spread after the debunking videos were posted. To some extent, these statistics reveal the universality and insufficient utilization of debunking videos.

Based on the above observations, we conjecture that the relationship among videos of the same event can be modeled to enhance the fake news video detection and rectify the detection results via factual information. To this end, we introduce the new cross-sample paradigm for fake news video detection and propose a corresponding novel framework, Nighbor-Enhanced fake Enews video Detection (**NEED**), which integrates the neighborhood relationship both explicitly and implicitly for better detection. NEED is a model-agnostic framework, which can easily incorporate various single-sample detectors to yield further improvement. Thus, we first obtain the feature representa-

tion from pre-trained single-sample detectors and then refine the representation and final prediction with relationship modeling.

To compensate for the insufficient information in a single video, we organize the news videos in the same event in the form of graph to aggregate the neighborhood information (**Graph Aggregation**). Specifically, we leverage the attention mechanism on the event graph (Velickovic et al., 2018) to model the correlations between different nodes and dynamically aggregate these features. Furthermore, as mentioned before, there exists explicit relation between debunking and fake news videos, *i.e.*, refutations. Consequently, debunking videos can be adopted to rectify the false negative predictions, spotting the “hidden” fake news videos (**Debunking Rectification**). Specifically, we formulate a new inference task to discriminate whether the given debunking video can refute the given candidate video. For a given video pair, the refutations commonly exist in the textual descriptions of the same visual scenes, which inspires us to detect the textual conflict of the same visual representation. To fulfill the discrimination, we take the visual representations from the video copy detector to obtain visual consistency, and fuse it with the textual feature from the textual conflict detector via the attention mechanism. Then the fusion feature is used to classify the refutation relationship between the debunking and candidate videos. Given the proposed graph aggregation and debunking rectification modules, NEED can significantly improve the performance of base single-sample detectors trained with single-modal or multimodal data.

Our contributions are summarized as follows:

- We propose a new cross-sample paradigm for fake news video detection, modeling multiple news videos in the same event simultaneously. Derived from such a paradigm, we propose the NEED framework, which exploits the neighborhood relationship explicitly and implicitly to enhance the fake news video detection.
- To the best of our knowledge, we are the first to utilize debunking videos in fake news video detection, which can utilize factual information to rectify false negative predictions. To this end, we formulate a new multimodal inference task and propose a novel model that utilizes the consistency from both the textual and visual perspectives to identify whether the given debunking video can refute the given

Debunking videos are videos that use factual evidence to refute widely circulated fake news, usually posted by experts.

candidate video.

- NEED is versatile and can be applied to various single-sample detectors. Extensive experiments on the public benchmark demonstrate that NEED can yield significant improvement with both single-modal and multimodal base detectors.

2 Related Work

To defend against fake news, researchers are mainly devoted to two threads of techniques:

Fake news detection methods commonly use non-factual multimodal signals such as linguistic patterns (Przybyla, 2020), image quality (Qi et al., 2019; Cao et al., 2020), multimodal inconsistency (Zhou et al., 2020; Qi et al., 2021), user response (Shu et al., 2019), and propagation structure (Ma et al., 2017), to classify the given news post as real or fake. With the prevalence of short video platforms, detecting fake news videos draws more attention in the community. Recent works mainly leverage deep neural networks to extract the multimodal features and model the cross-modal correlations (Choi and Ko, 2021; Shang et al., 2021; Palod et al., 2019; Qi et al., 2023). For example, Qi et al. (2023) use the cross-attention transformer to fuse news content features of different modalities including text, keyframes, and audio, and use the self-attention transformer to fuse them with social context features including comments and user.

However, existing works in fake news video detection identify each target news independently, without considering the neighborhood relationship in an event. In view of the practicality of the event-level process, Wu et al. (2022) construct a cross-document knowledge graph and employ a heterogeneous graph neural network to detect misinformation. Nonetheless, this work is performed on the synthetic dataset where each fake news document originates from a manipulated knowledge graph, which cannot be readily applied to real-world scenarios with unpredictable noises in information extraction. Moreover, they only consider the implicit relation among news texts while ignoring the explicit refutations between debunking information and fake news.

Fact-checking methods commonly rely on retrieved relevant factual information from reliable sources such as Wikipedia (Thorne et al., 2018) and webpages (Nie et al., 2019) to judge the veracity of the given check-worthy claim (Guo et al.,

2022; Zeng et al., 2021). A recent thread is to determine whether a claim has been previously fact-checked before retrieving evidence (Sheng et al., 2021). This task is commonly framed as a ranking task, ranking fact-checking articles based on the similarities to the given claim. Compared to textual fact-checking, multimodal verification is under-explored. Mishra et al. (2022) treat the verification as a multimodal entailment task, where the model needs to classify the relationship between the given reliable document (text with associated image) and check-worthy claim (text with associated image). Inspired by these works, the debunking rectification module in NEED focuses on rectifying the wrong predictions of previously fact-checked news videos by identifying the refutation relationship between the given debunking and candidate news video.

In summary, fake new detection methods leverage non-factual patterns learned from large-scale data to give timely judgments for newly emerging events, while fact-checking techniques provide more reliable judgments benefiting from the factual information but only work for a part of events limited by the coverage of external sources. Our work combines the merits of these two approaches: (1) We leverage the data-driven fake news video detectors to obtain effective multimodal representations and to model the neighborhood information, and (2) we also embrace the concept of relevant factual information in fact-checking to rectify the detection results with reliable debunking videos.

3 Methodology

3.1 Overview

As mentioned in the Introduction, the fabrication and malicious modification of fake news videos limit the verification ability of existing single-sample fake news video detectors, leading to inferior performance. In contrast, the relationship among neighborhood videos, *i.e.*, videos of the same event, can be used to supplement the current techniques. Thus, we propose the Neighbor-Enhanced fake News Video Detection (NEED) framework, leveraging the set of videos in an event, including fake news I_F , real news I_R and debunking videos I_D , to improve the performance of single-sample detectors. Specifically, NEED is model-agnostic, which takes the representations from the pre-trained base detectors (Feature Extraction) to build the dynamic graph and aggregate neighborhood information (Graph Aggregation, GA). Then,

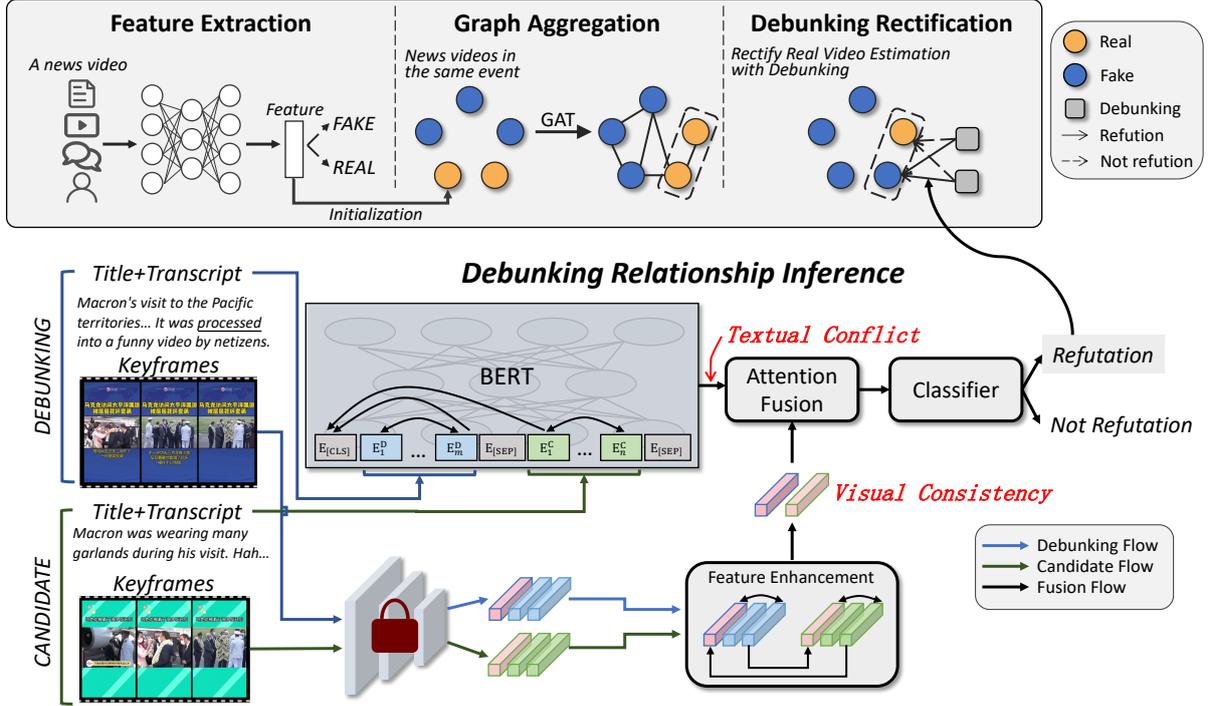


Figure 2: Architecture of the proposed framework NEED. The first row indicates the three stages in NEED, including feature extraction, graph aggregation, and debunking rectification. To realize the debunking rectification, debunking relation inference (the second row) is introduced to determine the refutation relationship.

we use the factual information from debunking videos to rectify the predicted results (Debunking Rectification, DR). The overall framework is illustrated in Figure 2.

3.2 Feature Extraction

News videos contain multimodal information, including title, audio, keyframes, video clips, comments, user profile, *etc.* Existing single-sample fake news video detectors leverage single-modal (Medina Serrano et al., 2020) or multimodal (Qi et al., 2023) information to discriminate each news video independently. They commonly design tailor-made modules to extract and fuse multimodal features. In contrast, NEED is a solution for the cross-sample paradigm, which can incorporate various single-sample fake news video detectors to yield further improvement with the neighborhood modeling. Thus, we first extract single-modal/multimodal features F_{base} for the given set of news videos from the base single-sample detector.

3.3 Graph Aggregation

Graph Construction. Given the set of related news video features F_{base}^E under the same event E , we organize them in the form of graph attention networks (GAT) (Velickovic et al., 2018). \mathcal{G} denotes the graph, \mathcal{V} denotes nodes in \mathcal{G} and \mathcal{E}

denotes edges between nodes. Each node $v_i \in \mathcal{V}$ represents a news video feature from the base detector, and the edge e_{ij} indicates the importance of node j 's feature to that of node i , which is obtained via attention mechanism.

Feature Aggregation and Classification. To aggregate the neighbor information, we apply the attention mechanism on the constructed event graph \mathcal{G} to update the representations of nodes. Specifically, given a node v_i with its neighbors \mathcal{N}_i , the weight α_{ij} between v_i and its neighbor $v_j \in \mathcal{N}_i$ is formulated as:

$$e_{ij} = \text{LeakyReLU}(\mathbf{a}^\top [\mathbf{W}v_i, \mathbf{W}v_j]),$$

$$\alpha_{ij} = \text{softmax}_j(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}_i} \exp(e_{ik})}, \quad (1)$$

where \mathbf{a} and \mathbf{W} are trainable parameters, \top denotes the matrix transpose, and $[\cdot, \cdot]$ is the concatenation operation. Then, the embedding of v_i is updated by the aggregated information:

$$\hat{v}_i = \sigma\left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W}v_j\right), \quad (2)$$

where σ is the nonlinear operation. To avoid over-smoothing of node features, we only adopt two GAT layers. The final feature \hat{v}_i is fed into a binary classifier to verify the video. The network is

optimized by the binary cross-entropy loss:

$$\mathcal{L} = -[(1 - y) \log(1 - p_{\text{GA}}) + y \log p_{\text{GA}}], \quad (3)$$

where p_{GA} is the predicted probability and $y \in \{0, 1\}$ denotes the ground-truth label.

3.4 Debunking Rectification

Graph aggregation focuses on combining the neighborhood features obtained from base detectors, which learn non-factual patterns from large-scale data. Instead, there also exists an explicit relationship between fake news videos and debunking videos with factual information, *i.e.*, refutations. Thus, we design the debunking rectification module to rectify the false negative predictions in the previous stages.

Specifically, we propose a new multimodal inference task to recognize this relationship, *i.e.*, debunking relationship inference. The definition of this task is as follows:

Definition 1: *Given a debunking video and a candidate video that belong to the same event, debunking relationship inference (DRI) aims to determine whether the debunking video can refute the candidate video or not.*

For a given event, we regard videos that are detected to be real by the GA module as the candidates $I_C = \{\eta_C^1, \dots, \eta_C^{n_c}\}$. For each candidate video η_C^i , we feed it into the DRI model together with the debunking videos $I_D = \{\eta_D^1, \dots, \eta_D^{n_d}\}$ in the same event. Then the candidate video is verified by combing the predicted probabilities of graph aggregation p_{GA}^i and DRI model p_{DR}^i :

$$\begin{aligned} p^i &= \max\{p_{\text{GA}}^i, p_{\text{DR}}^i\}, \\ p_{\text{DR}}^i &= \max_{\eta_D^j \in I_D} \text{DRI}(\eta_C^i, \eta_D^j). \end{aligned} \quad (4)$$

To realize the aim of DRI, we design the model following three principles: 1) Detecting the conflict between the news text of the debunking and candidate videos. 2) Detecting the consistency between video clips of the given video pair. For example, if the debunking video refutes a piece of fake news that misuses the “old” video clip from a previous event, we need to distinguish whether the candidate video uses this “old” video clip. 3) Dynamically fusing the textual and visual evidence to eliminate the irrelevant visual information for news events where the visual evidence is not essential, such as

“UN announces Chinese as the international common language”.

Based on the above principles, we propose a novel DRI model, which can detect and dynamically fuse textual conflict and visual consistency.

Textual Conflict Detection. Inspired by the task of natural language inference (NLI) (Bowman et al., 2015), we detect the textual conflict via the consistency between the given sentence pair. Specifically, given the debunking video, we extract and concatenate the title and video transcript as $S_D = [w_1, \dots, w_m]$, where w_i represents the i -th word in the composed sentence. Likewise, the news text in the candidate news video is represented as $S_C = [w_1, \dots, w_n]$. Then we pack the sentence pair $\langle S_D, S_C \rangle$ and feed it into BERT to model the intra- and inter- sentence correlations. The BERT we used has been fine-tuned on several NLI datasets to enhance its reasoning ability. A learnable type embedding is added to every token indicating whether it belongs to S_D or S_C . Finally, we obtain the textual conflict feature:

$$x_t = \text{BERT}([\text{CLS}]S_D[\text{SEP}]S_C[\text{SEP}]). \quad (5)$$

Visual Consistency Evaluation. To match the video clips, we leverage the EfficientNet (Tan and Le, 2021) pre-trained on the image similarity dataset (Douze et al., 2021) to obtain visual representations of each keyframe. We denote the frame features of the given debunking video and candidate video as $F_D = [f_D^1, \dots, f_D^l]$ and $F_C = [f_C^1, \dots, f_C^k]$, respectively. Following He et al. (2023), the fixed sine and cosine temporal positional encoding f_{tem} are added to the initial features, and a learnable classification token $f^{[\text{CLS}]}$ is prepended to the feature sequence as the global feature. The processed features of debunking video \hat{F}_D and candidate video \hat{F}_C are presented as:

$$\begin{aligned} \hat{F}_D &= [f_D^{[\text{CLS}]}, f_D^1, \dots, f_D^l] + f_{\text{tem}}, \\ \hat{F}_C &= [f_C^{[\text{CLS}]}, f_C^1, \dots, f_C^k] + f_{\text{tem}}. \end{aligned} \quad (6)$$

Similar to textual conflict detection, we need to consider intra- and inter- video correlations. Therefore, we employ stacked self- and cross- attention (Vaswani et al., 2017) modules to enhance the initial features, where the query vectors are from the other video in the cross-attention module. Finally,

As clarified in Fujian Province Debunking (2022), the fact is that there is no such thing as “international common language.”

the visual consistency feature is obtained by concatenating the classification tokens of the debunking and candidate videos:

$$\mathbf{x}_v = [\mathbf{f}_D^{[CLS]}, \mathbf{f}_C^{[CLS]}]. \quad (7)$$

Attention Fusion and Classification. Given the textual conflict feature \mathbf{x}_t and the visual consistency feature \mathbf{x}_v , we dynamically fuse them to spot the important information and eliminate irrelevant information via a self-attention fusion layer. Finally, the fused feature is fed into a binary classifier to estimate the probability p_{DR}^i in Eq. 4 that the debunking video can refute the candidate video.

4 Experiments

In this section, we conduct experiments to evaluate the effectiveness of NEED. Specifically, we aim to answer the following evaluation questions:

- **EQ1:** Can NEED improve the performance of fake news video detection?
- **EQ2:** How effective are the different modules of NEED in detecting fake news videos?
- **EQ3:** How does NEED perform in early detection, which means the number of videos in each event is limited?
- **EQ4:** How does NEED perform in the temporal split?

4.1 Experimental Setup

Dataset. We conducted experiments on the FakeSV dataset (Qi et al., 2023), the only fake news video dataset that provides rich events and debunking samples. This dataset collects news videos from popular Chinese short video platforms such as Douyin (the equivalent of TikTok in China), and employs human annotations. FakeSV consists of 1,827 fake news videos, 1,827 real news videos, and 1,884 debunked videos under 738 events. For each news video, this dataset provides the video, title, metadata, comments and user profile. Table 1 shows the statistics of this dataset.

Table 1: Statistics on the number of news videos in each event.

| | #Fake | #Real | #Debunking | All |
|------|-------|-------|------------|-----|
| Avg. | 3 | 3 | 3 | 8 |
| Min. | 0 | 0 | 0 | 1 |
| Max. | 24 | 21 | 20 | 25 |

douyin.com

Evaluation Metrics. To mitigate the performance bias caused by the randomness of data split, we follow the setting in Qi et al. (2023) and conduct evaluations by doing five-fold cross-validation with accuracy (Acc.), macro precision (Prec.), macro recall (Recall), and macro F1-score (F1) as evaluation metrics. For each fold, the dataset is split at the event level into a training set and a testing set with a sample ratio of 4:1. This ensures that there is no event overlap between different sets, thus avoiding the model detecting fake news videos by memorizing the event information (Wang et al., 2018).

Implementation Details. We use two GAT layers in GA and set the hidden states as 128 and 2, respectively, with ReLU for the first GAT layer. To avoid overfitting, a dropout layer is added between the two layers with a rate of 0.3. In DR, we use the pre-trained Erlangshen-MegatronBert-1.3B-NLI to evaluate the textual conflict. For visual consistency evaluation, we use the pre-trained EfficientNet to extract the frame features and use the pre-trained weight in the feature enhancement module. To train the debunking relationship inference model, the debunking videos and fake news videos in the same event are paired with the label “refutation”, and the debunking videos and real news videos are paired with the label “not refutation”. In the attention fusion module, we use a 4-head transformer layer. The last two layers of BERT, the visual module and the attention fusion module are trained for 30 epochs with a batch size of 64. The learning rate is set as 1×10^{-3} and 5×10^{-5} in GA and DRI, respectively. All experiments were conducted on NVIDIA RTX A5000 GPUs with PyTorch.

4.2 Base Models

NEED can readily incorporate any fake news video detectors that can produce video representation. Here we select four representative single-modal methods and two multimodal methods used in fake news video detection as our base detectors.

Single-modal: 1) **BERT** (Devlin et al., 2019) is one of the most popular textual encoders in NLP-related works. We concatenate the video caption and video transcript as a sequence and feed it into BERT for classification. 2) **Faster R-CNN+Attention** (Ren et al., 2015; Vaswani et al., 2017) is widely used in existing works (Shang et al.,

<https://github.com/IDEA-CCNL/Fengshenbang-LM>
<https://github.com/lyakaap/ISC21-Descriptor-Track-1st>
<https://github.com/transvcl/TransVCL>

Table 2: Performance (%) comparison of base models with and without NEED. The better result in each group using the same base model are in **boldface**, and the absolute gain is calculated. We report the mean and standard deviation of the five-fold cross-validation.

| | Method | Acc. | | F1 | | Prec. | | Recall | |
|--------------|-------------------|-------------------------|-----------------|-------------------------|-----------------|-------------------------|-----------------|-------------------------|-----------------|
| single-modal | BERT | 77.05 \pm 3.24 | – | 77.02 \pm 3.27 | – | 77.21 \pm 3.12 | – | 77.07 \pm 3.20 | – |
| | + NEED | 82.99 \pm 3.86 | 5.94 \uparrow | 82.96 \pm 3.87 | 5.94 \uparrow | 83.19 \pm 3.87 | 5.98 \uparrow | 82.99 \pm 3.88 | 5.92 \uparrow |
| | Faster R-CNN +Att | 70.19 \pm 2.70 | – | 70.00 \pm 2.68 | – | 70.68 \pm 2.89 | – | 70.15 \pm 2.69 | – |
| | + NEED | 78.48 \pm 3.30 | 8.29 \uparrow | 78.45 \pm 3.28 | 8.45 \uparrow | 78.71 \pm 3.45 | 8.03 \uparrow | 78.50 \pm 3.28 | 8.35 \uparrow |
| | VGGish | 66.91 \pm 1.33 | – | 66.82 \pm 1.30 | – | 67.07 \pm 1.41 | – | 66.89 \pm 1.32 | – |
| | + NEED | 75.25 \pm 1.61 | 8.34 \uparrow | 75.12 \pm 1.63 | 8.30 \uparrow | 75.73 \pm 1.67 | 8.66 \uparrow | 75.22 \pm 1.61 | 8.33 \uparrow |
| Multimodal | Wu et al. (2022) | 77.10 \pm 2.04 | – | 74.71 \pm 2.13 | – | 76.43 \pm 2.16 | – | 73.98 \pm 2.05 | – |
| | + NEED | 82.96 \pm 3.42 | 5.86 \uparrow | 82.93 \pm 3.44 | 8.22 \uparrow | 83.14 \pm 3.44 | 6.71 \uparrow | 82.95 \pm 3.46 | 8.97 \uparrow |
| | FANVM | 76.00 \pm 2.29 | – | 75.98 \pm 2.30 | – | 76.07 \pm 2.28 | – | 76.01 \pm 2.30 | – |
| | + NEED | 80.97 \pm 4.05 | 4.97 \uparrow | 80.90 \pm 4.10 | 4.92 \uparrow | 81.36 \pm 3.96 | 5.29 \uparrow | 80.96 \pm 4.04 | 4.95 \uparrow |
| | SV-FEND | 79.95 \pm 1.97 | – | 79.89 \pm 2.01 | – | 80.23 \pm 1.78 | – | 79.94 \pm 1.98 | – |
| | + NEED | 84.62 \pm 2.13 | 4.67 \uparrow | 84.61 \pm 2.12 | 4.72 \uparrow | 84.81 \pm 2.24 | 4.58 \uparrow | 84.64 \pm 2.14 | 4.70 \uparrow |

2021; Qi et al., 2023) to extract and fuse the visual features of multiple frames for classification. 3) **VGGish** (Hershey et al., 2017) is used to extract the acoustic features for classification. 4) **Wu et al. (2022)** construct a cross-document textual knowledge graph and employ a heterogeneous graph neural network to detect, which is one of the few works considering the cross-document relationship in fake news detection.

Multimodal: 1) **FANVM** (Choi and Ko, 2021) use topic distribution differences between the video title and comments as fusion guidance, and concatenate them with keyframe features. An adversarial neural network is used as an auxiliary task to help extract topic-agnostic multimodal features. 2) **SV-FEND** (Qi et al., 2023) use two cross-modal transformers to model the mutual enhancement between text and other modalities (*i.e.*, audio and keyframes), and then fuse them with social context features (*i.e.*, comments and user) by self-attention mechanism. Both of these multimodal methods are tailor-made for fake news video detection.

4.3 Performance Comparison (EQ1)

We compare the performance of base models with and without NEED in Table 2 and make the following observations: 1) With the help of NEED, all six base models gain significant performance improvement (4.67 ~ 8.34% in terms of accuracy), which validates the effectiveness and versatility of NEED. 2) Compared with **Wu et al. (2022)** that combines cross-document information, its basic feature encoder enhanced by NEED (*i.e.*, BERT+NEED) achieves better performance, verifying the superiority of NEED in utilizing the neighborhood correlations. 3) NEED yields more significant improvement

Table 3: Ablation studies on each component in NEED. GA: graph aggregation, DR: debunking rectification. The standard deviation values are ignored for simplicity.

| Method | Acc. | F1 | Prec. | Recall |
|----------------|--------------|--------------|--------------|--------------|
| SV-FEND | 79.95 | 79.89 | 80.23 | 79.94 |
| + DR | 80.94 | 80.90 | 81.15 | 80.93 |
| + GA | 83.43 | 83.41 | 83.61 | 83.45 |
| + NEED (DR&GA) | 84.62 | 84.61 | 84.81 | 84.64 |
| VGGish | 66.91 | 66.82 | 67.07 | 66.89 |
| + DR | 72.84 | 72.70 | 73.30 | 72.84 |
| + GA | 74.83 | 74.64 | 75.54 | 74.80 |
| + NEED (DR&GA) | 75.25 | 75.12 | 75.73 | 75.22 |
| DR | 82.95 | 81.05 | 81.36 | 81.04 |

on the underperformed model, *e.g.*, 8.34% improvement in Acc. on VGGish. We conjecture that such a phenomenon can be contributed to the explicit neighborhood modeling in the debunking rectification module, which ensures the lower bound of detection performance via factual information.

4.4 Ablation Studies (EQ2)

To verify the effectiveness of each proposed component in NEED, we conduct ablation experiments on top of both SOTA (*i.e.*, SV-FEND (Qi et al., 2023)) and underperformed (*i.e.*, VGGish (Hershey et al., 2017)) models in Table 2. From Table 3, we see that DR and GA consistently improve the performance of both base detectors. Moreover, comparing the two enhanced models, DR is more effective on the underperformed model than the SOTA model, which supports the explanation that DR ensures the lower bound of detection performance.

Interestingly, the improvement of DR is less significant than GA, especially on the SOTA model. We conjecture the reason lying in the limited de-

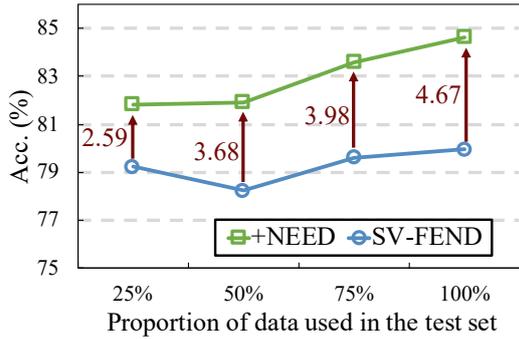


Figure 3: Performance of NEED in early detection.

Table 4: Performance of NEED under the temporal split.

| Method | Acc. | F1 | Prec. | Recall |
|---------|--------------|--------------|--------------|--------------|
| SV-FEND | 82.20 | 81.47 | 82.89 | 80.99 |
| +NEED | 89.67 | 89.37 | 90.16 | 88.97 |

bunking videos, which are only available in 51% of events in the FakeSV dataset. To further verify the effectiveness of factual information introduced by DR, we experiment with DR on the subset that contains debunking videos. Specifically, p_{DR}^i in Eq. 4 is used as the probability that the candidate video is fake. As shown in the last row in Table 3, solely using DR can achieve an accuracy of 82.95% on the subset, verifying the strong discriminability of debunking videos in detecting fake news videos. All the above results demonstrate that the neighborhood relationship can enhance and rectify fake news video detection.

4.5 Practical Settings

Early Detection (EQ3). Detecting fake news in its early stage is important for timely mitigating its negative influences (Guo et al., 2021). In this part, we conduct experiments using different data proportions of the test set to evaluate the performance of NEED with limited neighbors. Specifically, we keep the first 25%, 50%, 75% and 100% videos in each test event in chronological order, and conduct experiments on top of the SOTA base model SV-FEND. Figure 3 shows that NEED improves the base model even though with limited neighbors. Furthermore, as the number of videos within an event increases, NEED yields more significant improvement (from 2.59% at 25% data to 4.67% at 100% data), benefiting from the richer neighborhood relationship.



Figure 4: Illustration of the effect of graph aggregation. The left one is the video used in this event, and the right graph shows the score transformation of fake news videos before and after using GA.

Performance in Temporal Split (EQ4). Splitting data at the event level helps models learn event-invariant features and thus benefit generalization on new events, which is a common practice in the community (Wang et al., 2018; Qi et al., 2021). But in real-world scenarios, when a check-worthy news video emerges, we only have the previously-emerging data to train the detector. Thus we provide another temporal data split, which means splitting the dataset into training, validation and testing sets with a ratio of 70%:15%:15% in chronological order, to evaluate the ability of models to detect future fake news videos. Table 4 shows the performance of SV-FEND with and without NEED in the temporal split. We can see that NEED significantly improves the base model by 7.47% in Acc., demonstrating that the neighborhood relationship learned by NEED can readily benefit the detection of future fake news videos.

4.6 Case Studies

In this part, we list some cases to intuitively illustrate the effect of GA and DR.

Graph Aggregation Compensates Single Video Information. A single news video contains limited information, and the representation from single-sample detectors can be biased to some data patterns, such as verified publishers. Figure 4 shows the score transformation of multiple fake news videos in the same event before and after using GA. We infer that GA helps by transferring the key clue, *i.e.*, the indicative comment, in a single video to others. Moreover, by combining the neighbor information, GA mitigates the publisher bias of single-sample detectors (*i.e.*, videos published by verified users are commonly considered to be real).

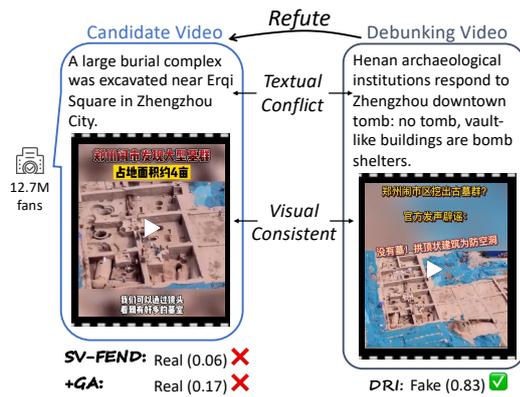


Figure 5: An example where the debunking video helps spot the “hard” fake news video missed by previous modules. The number denotes the predicted probability of labels being 0 (real) and 1 (fake), respectively.

Debunking Rectification Refutes Candidates via Factual Evidence. As shown in Figure 5, despite aggregating neighbor information ameliorates the biased prediction (probability 0.06 \rightarrow 0.17) based on the powerful publisher (verified institutional account with 12.7M fans), GA fails to address such a hard case with a strong bias. Instead, DR uses the debunking video with factual evidence to refute the candidate video, which successfully rectifies the false negative prediction.

5 Conclusion

We proposed a novel framework, namely NEED, to utilize the neighborhood relationship in the same event for fake news video detection. We designed the graph aggregation and debunking rectification modules to assist existing single-sample fake news video detectors. Experiments show the effectiveness of NEED in boosting the performance of existing models. We also drew insights on how the graph aggregation and debunking rectification contribute to fake news video detection.

Limitations

This work requires that news videos are organized into different events and each event has more than one candidate video. The debunking rectification module relies on the existence of labeled debunking videos, and the graph aggregation module relies on existing fake news detectors to provide the initial features for each video. The textual length in videos is limited due to that the debunking inference module is based on a pre-trained BERT model with limited sequence length.

Ethics Statement

Our framework in general does not create direct societal consequences and is intended to be used to defend against fake news videos. It can be easily combined into fake news video detection systems, especially when the events have multiple related news videos and debunking videos. To the best of our knowledge, no code of ethics was violated throughout the experiments done in this article. Experiments are conducted on the publicly available dataset and have no issues with user privacy.

Acknowledgements

This work was supported by the National Natural Science Foundation of China (62203425), the Zhejiang Provincial Key Research and Development Program of China (No.2021C01164), the Project of Chinese Academy of Sciences (E141020), the Innovation Funding from the Institute of Computing Technology, the Chinese Academy of Sciences under (E161020).

References

- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. *A large annotated corpus for learning natural language inference*. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 632–642. The Association for Computational Linguistics.
- Brandy Zadrozny. *On tiktok, audio gives new virality to misinformation* [online]. 2021.
- Juan Cao, Peng Qi, Qiang Sheng, Tianyun Yang, Junbo Guo, and Jintao Li. 2020. *Exploring the role of visual content in fake news detection*. *Disinformation, Misinformation, and Fake News in Social Media: Emerging Research Challenges and Opportunities*, pages 141–161.
- Hyewon Choi and Youngjoong Ko. 2021. *Using topic modeling and adversarial neural networks for fake news video detection*. In *CIKM '21: The 30th ACM International Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021*, pages 2950–2954. ACM.
- Cyberspace Administration of China. *The cyberspace administration of china guides the website platform to strengthen the traceability and disposal of online rumors related to the crash of the china eastern airlines crash incident* [online]. 2022. in Chinese.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.
- Matthijs Douze, Giorgos Tolias, Ed Pizzi, Zoë Papanikou, Lowik Chanussot, Filip Radenovic, Tomáš Jeníček, Maxim Maximov, Laura Leal-Taixé, Ismail Elezi, Ondrej Chum, and Cristian Canton-Ferrer. 2021. [The 2021 image similarity dataset and challenge](#). *CoRR*, abs/2106.09672.
- Dina ElBoghdady. 2013. [Market quavers after fake ap tweet says obama was hurt in white house explosions](#). *The Washington Post*.
- Marc Fisher, John Woodrow Cox, and Peter Hermann. 2016. [Pizzagate: From rumor, to hashtag, to gunfire in dc](#). *The Washington Post*, 6:8410–8415.
- Fujian Province Debunking. [Un announces chinese as the international common language? fake!](#) [online]. 2022.
- Bin Guo, Yasan Ding, Lina Yao, Yunji Liang, and Zhiwen Yu. 2021. [The future of false information detection on social media: New perspectives and trends](#). *ACM Comput. Surv.*, 53(4):68:1–68:36.
- Zhijiang Guo, Michael Sejr Schlichtkrull, and Andreas Vlachos. 2022. [A survey on automated fact-checking](#). *Trans. Assoc. Comput. Linguistics*, 10:178–206.
- Sifeng He, He Yue, Minlong Lu, et al. 2023. [Transvcl: Attention-enhanced video copy localization network with flexible supervision](#). In *37th AAAI Conference on Artificial Intelligence: AAAI 2023*.
- Shawn Hershey, Sourish Chaudhuri, Daniel P. W. Ellis, Jort F. Gemmeke, Aren Jansen, R. Channing Moore, Manoj Plakal, Devin Platt, Rif A. Saurous, Bryan Seybold, Malcolm Slaney, Ron J. Weiss, and Kevin W. Wilson. 2017. [CNN architectures for large-scale audio classification](#). In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2017, New Orleans, LA, USA, March 5-9, 2017*, pages 131–135. IEEE.
- Rui Hou, Verónica Pérez-Rosas, Stacy L. Loeb, and Rada Mihalcea. 2019. [Towards automatic detection of misinformation in online medical videos](#). In *International Conference on Multimodal Interaction, ICMI 2019, Suzhou, China, October 14-18, 2019*, pages 235–243. ACM.
- Jing Ma, Wei Gao, and Kam-Fai Wong. 2017. [Detect rumors in microblog posts using propagation structure via kernel learning](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 708–717. Association for Computational Linguistics.
- Juan Carlos Medina Serrano, Orestis Papakyriakopoulos, and Simon Hegelich. 2020. [NLP-based feature extraction for the detection of COVID-19 misinformation videos on YouTube](#). In *Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020*, Online. Association for Computational Linguistics.
- Shreyash Mishra, Suryavardan S, Amrit Bhaskar, Parul Chopra, Aishwarya N. Reganti, Parth Patwa, Amitava Das, Tanmoy Chakraborty, Amit P. Sheth, and Asif Ekbal. 2022. [FACTIFY: A multi-modal fact verification dataset](#). In *Proceedings of the Workshop on Multi-Modal Fake News and Hate-Speech Detection (DE-FACTIFY 2022) co-located with the Thirty-Sixth AAAI Conference on Artificial Intelligence (AAAI 2022), Virtual Event, Vancouver, Canada, February 27, 2022*, volume 3199 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Salman Bin Naeem and Rubina Bhatti. 2020. [The covid-19 ‘infodemic’: a new front for information professionals](#). *Health Information & Libraries Journal*, 37(3):233–239.
- Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. [Combining fact extraction and verification with neural semantic matching networks](#). In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 6859–6866. AAAI Press.
- Priyank Palod, Ayush Patwari, Sudhanshu Bahety, Saurabh Bagchi, and Pawan Goyal. 2019. [Misleading metadata detection on youtube](#). In *Advances in Information Retrieval - 41st European Conference on IR Research, ECIR 2019, Cologne, Germany, April 14-18, 2019, Proceedings, Part II*, volume 11438 of *Lecture Notes in Computer Science*, pages 140–147. Springer.
- Piotr Przybyla. 2020. [Capturing the style of fake news](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 490–497. AAAI Press.
- Peng Qi, Yuyan Bu, Juan Cao, Wei Ji, Ruihao Shui, Junbin Xiao, Danding Wang, and Tat-Seng Chua. 2023. [FakeSV: A multimodal benchmark with rich social context for fake news detection on short video platforms](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Peng Qi, Juan Cao, Xirong Li, Huan Liu, Qiang Sheng, Xiaoyue Mi, Qin He, Yongbiao Lv, Chenyang Guo,

- and Yingchao Yu. 2021. [Improving fake news detection by using an entity-enhanced framework to fuse diverse multimodal clues](#). In *MM '21: ACM Multimedia Conference, Virtual Event, China, October 20 - 24, 2021*, pages 1212–1220. ACM.
- Peng Qi, Juan Cao, Tianyun Yang, Junbo Guo, and Jintao Li. 2019. [Exploiting multi-domain visual information for fake news detection](#). In *2019 IEEE International Conference on Data Mining, ICDM 2019, Beijing, China, November 8-11, 2019*, pages 518–527. IEEE.
- Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. 2015. [Faster R-CNN: towards real-time object detection with region proposal networks](#). In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 91–99.
- Lanyu Shang, Ziyi Kou, Yang Zhang, and Dong Wang. 2021. [A multimodal misinformation detector for COVID-19 short videos on tiktok](#). In *2021 IEEE International Conference on Big Data (Big Data), Orlando, FL, USA, December 15-18, 2021*, pages 899–908. IEEE.
- Qiang Sheng, Juan Cao, Xueyao Zhang, Xirong Li, and Lei Zhong. 2021. [Article reranking by memory-enhanced key sentence matching for detecting previously fact-checked claims](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 5468–5481. Association for Computational Linguistics.
- Kai Shu, Limeng Cui, Suhang Wang, Dongwon Lee, and Huan Liu. 2019. [defend: Explainable fake news detection](#). In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019*, pages 395–405. ACM.
- Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. [Fake news detection on social media: A data mining perspective](#). *SIGKDD Explor.*, 19(1):22–36.
- S Shyam Sundar, Maria D Molina, and Eugene Cho. 2021. [Seeing is believing: Is video modality more powerful in spreading fake news via online messaging apps?](#) *Journal of Computer-Mediated Communication*, 26(6):301–319.
- Mingxing Tan and Quoc V. Le. 2021. [Efficientnetv2: Smaller models and faster training](#). In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 10096–10106. PMLR.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. [FEVER: a large-scale dataset for fact extraction and verification](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 809–819. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. [Graph attention networks](#). In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- 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 *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, pages 849–857. ACM.
- Xueqing Wu, Kung-Hsiang Huang, Yi R. Fung, and Heng Ji. 2022. [Cross-document misinformation detection based on event graph reasoning](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 543–558. Association for Computational Linguistics.
- Xia Zeng, Amani S. Abumansour, and Arkaitz Zubiaga. 2021. [Automated fact-checking: A survey](#). *Lang. Linguistics Compass*, 15(10).
- Xinyi Zhou, Jindi Wu, and Reza Zafarani. 2020. [SAFE: similarity-aware multi-modal fake news detection](#). In *Advances in Knowledge Discovery and Data Mining - 24th Pacific-Asia Conference, PAKDD 2020, Singapore, May 11-14, 2020, Proceedings, Part II*, volume 12085 of *Lecture Notes in Computer Science*, pages 354–367. Springer.

ACL 2023 Responsible NLP Checklist

A For every submission:

- A1. Did you describe the limitations of your work?
Limitations
- A2. Did you discuss any potential risks of your work?
Not applicable. Left blank.
- A3. Do the abstract and introduction summarize the paper's main claims?
abstract; Section 1
- A4. Have you used AI writing assistants when working on this paper?
Left blank.

B Did you use or create scientific artifacts?

Section 4

- B1. Did you cite the creators of artifacts you used?
Section 4
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
We were unable to find the license for the dataset we used.
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
Ethics Statement
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
Ethics Statement
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
Section 4
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
Section 4

C Did you run computational experiments?

Section 4, Appendix

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
Section 4

The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

Section 4

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

Section 4

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

Section 4

D Did you use human annotators (e.g., crowdworkers) or research with human participants?

Left blank.

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

No response.

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

No response.

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

No response.

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

No response.

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

No response.