# Multimodal Fake News Detection via CLIP-Guided Learning

Yangming Zhou
School of Computer Science
Shanghai, China
ymzhou21@fudan.edu.cn

Qichao Ying
School of Computer Science
Shanghai, China
qcying20@fudan.edu.cn

Zhenxing Qian*
School of Computer Science
Shanghai, China
zxqian@fudan.edu.cn

Sheng Li
School of Computer Science
Shanghai, China
lisheng@fudan.edu.cn

Xinpeng Zhang
School of Computer Science
Shanghai, China
zhangxinpeng@fudan.edu.cn

## ABSTRACT

Multimodal fake news detection has attracted many research interests in social forensics. Many existing approaches introduce tailored attention mechanisms to guide the fusion of unimodal features. However, how the similarity of these features is calculated and how it will affect the decision-making process in FND are still open questions. Besides, the potential of pretrained multimodal feature learning models in fake news detection has not been well exploited. This paper proposes a FND-CLIP framework, i.e., a multimodal Fake News Detection network based on Contrastive Language-Image Pretraining (CLIP). Given a targeted multimodal news, we extract the deep representations from the image and text using a ResNet-based encoder, a BERT-based encoder and two pairwise CLIP encoders. The multimodal feature is a concatenation of the CLIP-generated features weighted by the standardized cross-modal similarity of the two modalities. The extracted features are further processed for redundancy reduction before feeding them into the final classifier. We introduce a modality-wise attention module to adaptively reweight and aggregate the features. We have conducted extensive experiments on typical fake news datasets. The results indicate that the proposed framework has a better capability in mining crucial features for fake news detection. The proposed FND-CLIP can achieve better performances than previous works, i.e., **0.7%**, **6.8%** and **1.3%** improvements in overall accuracy on Weibo, Politifact and Gossipcop, respectively. Besides, we justify that CLIP-based learning can allow better flexibility on multimodal feature selection.

## CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

## KEYWORDS

datasets, neural networks, gaze detection, text tagging

## 1 INTRODUCTION

Online social networks have largely replaced the conventional way of information communication represented by newspapers and magazines. People enjoy the convenience of online social media in seeking friends or sharing viewpoints. However, OSNs have also promoted the wide and rapid spreading of fake news [24, 33, 49]. Online news posts can be more easily manipulated compared to written materials. News forgery can take various forms, for example, replacing a critical object within a picture with another one, or making biased or even misleading comments on the picture. What's worse, the readers are susceptible to well-crafted fake news and further circulate them. In sum, fake news is likely to create panic and misdirect public opinion, which alters society in negative ways.

In the past decades, Fake News Detection (FND) has been the center of data-centric research for decades [1, 36, 45]. While manual observation towards all news and posts on the Internet is both expensive and time-consuming, automatic FND using machine learning is an efficient way to combat the widespread dissemination of fake news. FND has helped news readers identify bias and misinformation in news articles and therefore stop their spreading. Early works on fake news detection merely focused on text-only or image-only content analysis [5, 11]. A pretrained model is usually employed to verify the logical and semantic soundness of the input. Also, trivial clues such as grammatical errors or traces left by image manipulation might be taken into consideration. While unimodal FND schemes are effective, modern news and posts are usually with rich information of several modalities and these methods neglect their correlation. For some fake news, a real image can be combined with total rumors and correct words can be used to describe a tampered image. In that sense, multimodal feature analysis is required to offer complementary benefits to assist FND.

In recent years, there have already been a lot of works that aggregate multimodal features to detect anomalies in news and posts [9, 24, 42]. Besides fusing features from images and texts,

**I have seen the biggest Wang Ba!** According to Weibo, around 11 noon today, in the river below the new bridge of Nanxun, Huzhou City, it was found to be more than 400 kg of turtle. EssenceEssenceTen people caught it.

Prediction:        Fake
Similarity Score:  0.829
Attention Score:   0.415 | 0.118 | 0.467

Netizens broke the news, at 7 pm on November 21, unknown flying objects over Deyang! Suspected of space airship! Unknown flying objects will glow, with the words Deyang Chengdu 2000 kilometers.

Prediction:        Real
Similarity Score:  0.859
Attention Score:   0.415 | 0.118 | 0.467

After the Kunming cutting incident, the Malay plane lost contact, and the Changsha cut people. Today, the Guangzhou Shahe on the morning of 3.15 has another incident. In the afternoon, there was another cut.

Prediction:        Fake
Similarity Score:  0.341
Attention Score:   0.392 | 0.211 | 0.397

Refugees at 8 Budapest Railway Station tell their own story. This railway station has only 8 toilets, but there are hundreds of refugees.

Prediction:        Real
Similarity Score:  0.277
Attention Score:   0.406 | 0.159 | 0.435

**Figure 1: Examples of fake news detection using FND-CLIP.The three attention scores of each news are text score, image score, and fusion score respectively.**

comments, up-vote ratio and the spreading graph are mostly preferred by researchers to evaluate the truthfulness of a post. These additional modalities are interactive and change over time. Many previous works prefer using as much modalities as possible. However, interactive modalities are less dependable than images and texts, or static modalities. First, the absence of interactive modalities might be common. A typical example is that no clue can be left in historical behavior in news posted by newly registered users, and nor can it be left in comments or up-votes if we wish to reject fake news shortly after their submission. Second, interactive modalities are less stable and can be changed over time, therefore potentially resulting in varied forensics results. Therefore, we revisit current arts in FND with only static modalities, and find that though many algorithms design well-crafted networks for multimodal feature fusion [38, 42], the mechanisms are largely at a black-box level as to how multimodal features will influence the final decision. Some works try to address this issue by explicitly calculating correlation on generating fused features. For example,

Chen et al. [9] additionally train variational auto encoders (VAE) that first compress the images and texts and contrastively learns to minimize the Kullback-Leibler (KL) divergence for news with correct image-text pairs. The corresponding cross-modal ambiguity score is then used to reweight the multimodal features [9]. Dhruv et al. [24] propose MVAE that trains a decoder to reconstruct the original texts and low-level image features from the fused features. These methods have achieved decent performance in multimodal fake news detection.

However, there are still some issues for multimodal FND to be addressed. First, we find that how the similarity of features from different modalities is to be calculated and how it will affect the decision-making process in FND is still an open question. For [9], we are not sure how efficient the VAEs are so that the KL divergence will be small given matched image-text pairs. For MVAE, though the ability of reconstruction means that the fused features are able to contain more information, the necessity of these auxiliary tasks in the view of FND remains unknown. Besides, we find that more advanced multimodal learning paradigms and pretrained models are not properly applied in FND. For example, CLIP [34] is a multimodal model that combines knowledge of language concepts with semantic knowledge of images. It was trained on a variety of image-text pairs to predict the most relevant text snippet, given an image, and vice versa. CLIP, together with other advanced multimodal technologies can be beneficial in image-text feature fusing, yet their usages in FND still remain ill-posed.

This paper proposes FND-CLIP, a multimodal fake news detection network based on the pretrained Contrastive Language-Image Pretraining (CLIP) model. The CLIP-based learning for fake news detection is to address the issue of cross-modal ambiguity by explicitly measuring the correlation between texts and images of targeted posts, and to guide the feature fusing and decision-making stages. Specifically, we encode the image using a fine-tunable ResNet [16] encoder a pretrained CLIP image encoder. The text is encoded by a fine-tunable BERT [13] encoder as well as a CLIP text encoder. The unimodal features are generated by concatenating the CLIP-generated features with the fine-tunable counterparts. The fused features consist of the two CLIP outputs. We use three projection heads to individually process the unimodal and fused features, which shrinks their sizes in order to distill the most important features for FND. Besides, we calculate the cosine similarity on the CLIP outputs and standardize it as the cross-modal similarity score. The score reweights the fused feature, where we regulate that less information will be provided by the fused features if the image and text show low correlation. Furthermore, we introduce an attention layer that outputs three scores that adaptively measure the significance of these features in their contribution to fake news detection. The classifier finally processes the summarized features to distinguish fake news from real ones.

We have conducted extensive experiments on FND-CLIP on several typical datasets for FND, including a Chinese dataset named Weibo, and two English datasets named Politifact and Gossip. The results show that FND-CLIP achieves **0.7%, 6.8%** and **1.3%** performance improvement in overall accuracy on the three datasets. Besides, we justify that CLIP-based learning can allow better flexibility on multimodal feature selection. Figure 1 showcases four examples of fake news detection using FND-CLIP, where we see that

the attention score as well as cross-modal similarity vary among different news instances. FND-CLIP is able to pay less attention to the multimodal features when the similarity is low, therefore flexibly aggregating information according to the characteristics of the provided news.

The contributions of this paper are mainly three-folded, namely:

- We propose FND-CLIP, a multimodal fake news detection method with CLIP-based learning, where the CLIP pretrained model is used to measure the cross-modal similarity and guide the mapping and fusion of features.
- We propose a modality-wise attention mechanism to adaptively weight the text, image, and fused features. Given different news instances, we find that the model flexibly learns to pay more attention to useful information in unimodal or multimodal features.
- We have conducted comprehensive experiments on three famous datasets, where the results prove that CLIP-generated features can be important assists to the unimodal features. FND-CLIP outperforms state-of-the-art fake news detection methods.

## 2 RELATED WORKS

### 2.1 Unimodal Fake News Detection

Unimodal FND usually works on finding anomalies in either the text or the image of a post. These algorithms often follow the essence of human decision process. For images, Cao et al. [7] jointly study image forensics features, semantic features, statistical features, and context features for fake news detection. It suggests that typical methods for image manipulation detecion [8] are useful in unveiling traces for news tampering. Besides, semantic inconsistency regarding the common sense [26] as well as poor image quality [15] can be widely present in fake news. For texts, verifying the logical soundness is essential [14], also accompanied by finding clues such as grammatical errors, writing styles [31] or extracting rhetorical structure [11]. Besides, both linguistic and visual patterns can be highly dependent on specific events and corresponding domain knowledge. Therefore, Nan et al. [29] propose to employ domain gate to aggregate multiple representations extracted by mixture-of-experts, and it deals with multi-domain fake news propagation in the language modality.

Though these unimodal characteristics can be explored and they indeed play key roles in distinguishing fake news, the multimodal characteristics such as correlation and consistency are ignored, which potentially impair the overall performance of these unimodal schemes on multimodal news.

### 2.2 Multimodal Fake News Detection

In the past literature, many works have been done on mining useful representations from images and texts of the news for fake news detection. Earlier works design sophisticated yet black-box attention mechanisms for multimodal feature fusion [3, 5]. Many other works [9, 24, 42] propose to better align the extracted features from different modalities before sending them into the classifier. Wang et al. [42] propose EANN that further employs an auxiliary task of event classification to aid feature extraction. The event classification branch is designed to better disentangle the mined multimodal

features so that there are both event-specific information and event-agnostic information. Dhruv et al. [24] processe the image and text using unimodal feature extractors and further utilize a multimodal VAE to learn a shared representation from them. The sampled representation produced by the VAE is then sent to a decoder which tries to reconstruct the original texts and low-level image features. Besides the focus on network design, other works exploit more information from the datasets. For example, Qi et al. [32] claim that image feature extractors cannot well understand visual entities such as celebrities, landmarks, and texts within the images, and therefore propose to manually extract these kinds of information as linguistic assists. Zhang et al. [47] design a novel dual emotion feature descriptor to measure the emotional gap between the post and its comments and verify that dual emotion is distinctive between fake and real news. Chen et al. [9] use two VAEs to compress the images and texts and contrastively learn to minimize the Kullback-Leibler (KL) divergence for correctly matched image-text pairs. The resultant score is then used to reweight the multimodal features during feature fusing.

Though these methods have achieved decent performance in multimodal FND, there are still issues to be concerned. First, how to explicitly measure the correlation between images and texts within a post still remain unclear. Second, we see that little work in FND consider applying the recently emerged arts in multimodal learning, which motivates us to use the CLIP-based pretraining to further boost the performance.
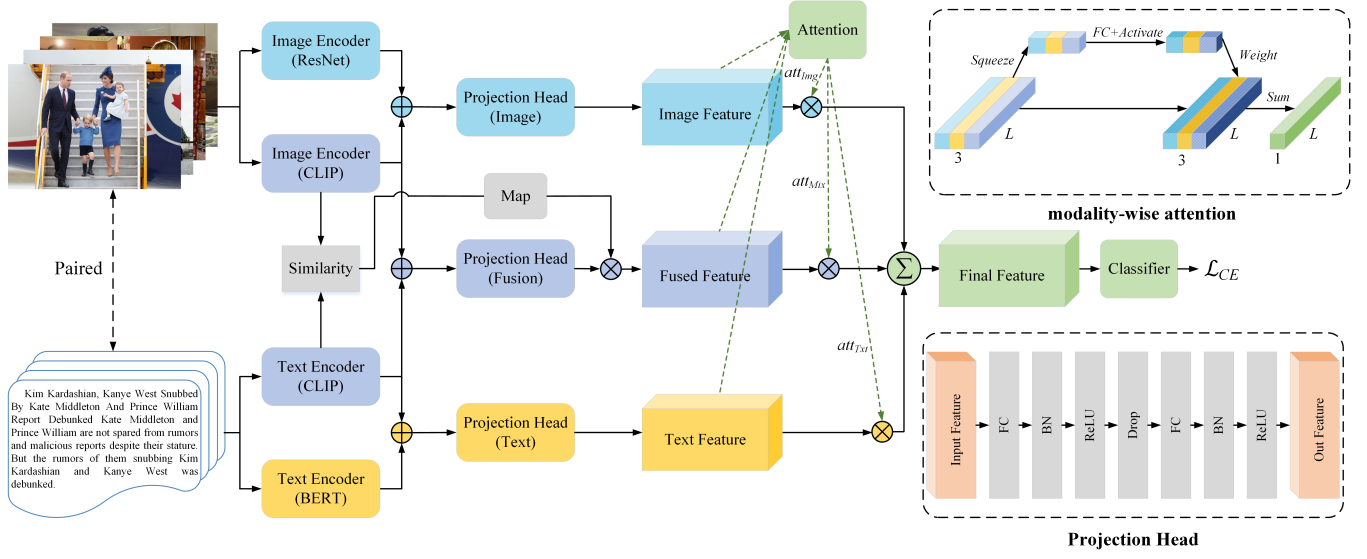
### 2.3 Multimodal Learning

Recent years have shown rapid developments in the field of multimodal machine learning [2]. Neural architectures are employed in tasks that go beyond single modalities, for example, Visual Question Answering (VQA) [12], Visual Commonsense Reasoning (VCR) [46], etc. In these tasks and beyond, priors and features from different modalities are required and algorithms or deep networks cannot be effective when provided with only a single modality. Several generic technologies are developed for learning joint representations of image content and natural language. For example, the CLIP model [34] is designed as a bridge between computer vision and natural language processing. It was trained on a variety of image-text pairs to predict the most relevant text snippet, given an image, without directly optimizing for the task. The model consists of two encoders that respectively embed texts and images into a uniform mathematical space. Then, for the matched image-text pair, CLIP is encouraged to maximize the cosine similarity between the embedding of the two modalities. Otherwise, the similarity is minimized for the model to find the most suitable paired images and texts. Multimodal learning has a promising future where the innovation of CLIP has benefitted a lot of down-stream tasks [10, 43]. Other multimodal schemes can be represented by Glide [30] and VilBERT [28] that are respectively for text-to-image generation and multimodal representation learning.

## 3 METHOD

### 3.1 Approach Overview

For multimodal fake news detection, we collect the statc modalities of the sampled news that includes text and image, and denote each

**Figure 2: The architecture of the proposed FND-CLIP method. CLIP, BERT, and ResNet are used to extract the features of different modalities of multimodal news. Encoded features of different levels are obtained through projection heads. CLIP similarity score is calculated to determine the importance of fused feature. A modality-wise attention mechanism is further used to reweight different modal features adaptively for the classifier to classify fake news.**

sample as $\mathbf{x} = (\mathbf{x}_{Txt}, \mathbf{x}_{Img})$. The ground-truth label is $y$ where $y = 0$ indicates that $\mathbf{x}$ is a real news, otherwise $y = 1$. According to the most traditional multimodal learning paradigm, a rich set of features are first extracted from $\mathbf{x}_{Txt}$ and $\mathbf{x}_{Img}$ that both represents the unimodal characteristics and the multimodal characteristics, which are then further fused and projected into a single value of $\hat{y}$ that should be close to the ground truth.

$$\hat{y} = F_{cls}(F_{Mix}(F_{Txt}(\mathbf{x}_{Txt}), F_{Img}(\mathbf{x}_{Img}))), \quad (1)$$

where $F_{Txt}$ and $F_{Img}$ are unimodal feature extractors, $F_{Mix}$ is the feature fusing model and $F_{cls}$ is the classification head. In order to model $F_{Txt}$ and $F_{Img}$, most of the previous methods use different pre-trained models to extract text and image features in different semantic spaces, and for $F_{Mix}$, the proposed mechanisms vary. The crucial point is how to ensure that features provided from both modalities will be utilized in the later stage, otherwise the gap in semantic space makes the fused features unable to accurately represent the correlation between image and text. Instead of applying sophisticated and black-box feature-fusing networks, we employ a simple yet effective method where pretrained networks for multimodal learning is introduced to extract aligned multimodal features and to guide the learning of the classification network. We choose the CLIP model [34] to measure the cross-modal similarity considering that the model is trained to provide the most appropriate language description of a given image and vice versa, and therefore is in line with the above requirements. After feature extraction and alignment, we use a light-weight network to implement $L_{Cls}$ which predicts $\hat{y}$.

## 3.2 Network Specification

Figure 2 illustrates the network design of FND-CLIP. The whole pipeline consists of four main modules, namely, unimodal feature encoder, CLIP-based encoder, projection and attention module, and finally the classifier.

**Unimodal feature generation.** We use a pretrained BERT model to obtain the feature $f_{BERT} \in \mathbb{R}^{n_{BERT}}$ of $\mathbf{x}_{Txt}$. For the image $\mathbf{x}_{Img}$, we use ResNet [17] to get deep representations $f_{ResNet} \in \mathbb{R}^{n_{ResNet}}$ from the image. Besides $f_{BERT}$ and $f_{ResNet}$, we use CLIP encoders to encode text and image and obtain the features $f_{CLIP\text{-}T} \in \mathbb{R}^{n_{CLIP}}$ and $f_{CLIP\text{-}I} \in \mathbb{R}^{n_{CLIP}}$. In order to improve the representation capability of the unimodal branches, embedding concatenation are performed in the text and image unimodal intra-modalities, respectively,

$$\begin{cases} f_{Txt} = concat(f_{BERT}, f_{CLIP\text{-}T}) \\ f_{Img} = concat(f_{ResNet}, f_{CLIP\text{-}I}), \end{cases} \quad (2)$$

where $f_{Txt} \in \mathbb{R}^{n_{BERT}+n_{CLIP}}$ and $f_{Img} \in \mathbb{R}^{n_{ResNet}+n_{CLIP}}$.

**CLIP-guide multimodal feature generation.** The text and image features extracted by BERT and ResNet respectively have significant cross-modal semantic gaps, and it is difficult for the network to learn their intrinsic semantic correlation if they are fused directly. Therefore, the two features are only used as unimodal representation, while the multimodal representation is obtained by first concatenating the alignment features of the text-image pair extracted by CLIP and then fine-tuning them to reduce redundancy and introduce attention. The concatenated feature is denoted as $f_{Mix} \in \mathbb{R}^{2 \times n_{CLIP}}$, where

$$f_{Mix} = concat(f_{CLIP-T}, f_{CLIP-I}). \quad (3)$$

The multimodal features reflect the correlation between the two modalities and contain meaningful semantic information. The assistance of the multimodal features to unimodal features is to learn the cross-modal similarity. Previous works often use a single network to mine both coarse and fine features from a modality, which is quite demanding on the learning ability of the model. Here, with the introduction of CLIP model, BERT and ResNet, which is the pre-training models for unimodal tasks, can pay more attention to trivial clues compared to extracting semantic information. For example, BERT can better extract emotional features of texts, and ResNet can identify higher-frequent noise patterns of images. In contrast, the training strategy of CLIP uses large-scale image-text pairs to learn the extraction of semantics, while largely ignoring emotion, noise and other features irrelevant to image and text matching. Therefore, using CLIP for multimodal feature generation can well collaborate with the unimodal features to respectively scrutinize the news from different aspects.

After we get the three features of different modalities, we use three individual projection head $P_{Txt}$, $P_{Img}$ and $P_{Mix}$ made up of Multi-Layer Perceptrons (MLP) to process the features. The goal is to reduce the dimension of the coarse features provided by the encoders and help filtering out redundant information. These networks share the same architecture but do not share weights. As is shown in Figure 2, every the projection head contains two sets of full connected layer with Batch Normalization [20] layer, a ReLU activation function, and a dropout layer.

Merely combining the CLIP-based features as the multimodal features cannot necessarily provide enough reliable information. The reason is that the authenticity of news is not completely correlated with image-text correlation. Some news posts, no matter real or fake, lack cross-modal relation or even semantic information. In that case, some instances require more emotion, noise, and other features, and the corresponding multimodal features might be noisy when the similarity is low and fully utilizing such information might impair the performance. To address the ambiguity issue between multimodal features, we measure the cosine similarity between the text features and the image features provided by CLIP, to adjust the intensity of fused features. The cosine similarity is calculated as follows.

$$sim = \frac{f_{Txt} \cdot (f_{Img})^T}{\|f_{Txt}\| \|f_{Img}\|}. \tag{4}$$

Then, we apply standardization and a Sigmoid functions to map the similarity into the range $[0 - 1]$. The normalization is done by calculating the running status of mean and standard deviation during training, and subtract the running mean from $sim$ and divide it with the running standard deviation. Compared to the contrastive learning paradigm, the normalization helps to calculate the similarity without comparing the news post with other instances.

Thus, the process of obtaining the projected unimodal and multimodal features is as follows.

$$\begin{cases} m_{Txt} = P_{Txt}(f_{Txt}) \\ m_{Img} = P_{Img}\left(f_{Img}\right) \\ m_{Mix} = Sigmoid\left(Std\left(sim\right)\right) \cdot P_{Mix}\left(f_{Mix}\right). \end{cases} \tag{5}$$

**Feature aggregation using modality-wise attention.** We apply an attention mechanism to reweight the projected features before aggregating the features from different modalities using spatial addition. Inspired by the Squeeze-and-Excitation Network (SE-Net) [19], we designed a modality-wise attention module as shown in Figure 2 to weight each feature adaptively. First, the three $L \times 1$ features are concatenated into one $L \times 3$ feature, where $L$ represents the length of the feature. Average pooling and maximum pooling are adopted to squeeze a $1 \times 3$ vector via summation, corresponding to the initial weight of each channel. Then, the initial weight obtained in the previous step is sent into the two $3 \times 3$ fully connected layers with GELU [18] activation function, and normalized into the range $[0 - 1]$ using Sigmoid functions respectively to obtain the attention weights $att = \{att_{Txt}, att_{Img}, att_{Mix}\}$. Finally, the weights are multiplied respectively on $m_t$, $m_i$ and $m_{mix}$, and a sum process is performed to obtain the $L \times 1$ aggregated feature $m_{Agg}$.

$$m_{Agg} = att_{Txt} \cdot m_{Txt} + att_{Img} \cdot m_{Img} + att_{Mix} \cdot m_{Mix}. \tag{6}$$

**Classification and objective function.** We feed the aggregated representation $m_{Agg}$ into a two-layer fully-connected network as the classifier $F_{cls}$ to predict the label $\hat{y}$. The objective function of FND-CLIP is to minimize the cross-entropy loss to correctly predict the real and fake news.

$$\mathcal{L}_{CE} = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}). \tag{7}$$

### 3.3 Training Detail

On the selection of BERT pretrained models, we respectively use the "bert-base-chinese" model on Chinese data and the "bert-base-uncased" model on English data, perform an attention-based post-processing [21]. The length of the input text is set to 300 words. About the ResNet, we use pre-trained ResNet-101 to extract visual features, setting the size of the input image to 224 × 224. The size of the images inputted to CLIP is the same as that to ResNet. Science CLIP has not pre-trained Chinese text model, we use Google Translation API [23] to translate Chinese texts to English. In addition, we use the summary generation model [35] to generate summary statements as the CLIP input for the text with the size longer than 50, to meet the requirements that the input size of the text has an upper bound in CLIP. The used pre-trained CLIP model is "ViT-B/32". We fine-tune ResNet in training stage, while freezing the weights of BERT and CLIP due to their difficulty in training on small datasets. We implement the projection heads using two fully connected layers with 256 and 64 hidden units, respectively. The hidden sizes of the two fully connected layers in the classifier are 64 and 2, respectively. The batch size is set as 64.

We use Adam optimizer [25] with the default parameters. The learning rate is $1 \times 10^{-3}$ where weight decay is 12. We trained a model for 50 epochs and chose the epoch getting the best test accuracy among them as the final result to avoid over-fitting.

## 4 EXPERIMENTS

### 4.1 Experimental Setup

**Dataset.** We use three real-world datasets collected from social media, namely, Weibo [22], Gossipcop, and Politifact [36]. During experiments, the unimodal news posts with no image or no text

**Table 1: Performance comparison between FND-CLIP and other methods on three datasets. Our method achieves the highest accuracy among these methods, and its precision, recall, and FI-score are also higher than most of the compared methods.**

| | Method | Accuracy | Fake News | | | Real News | | |
|---|---|---|---|---|---|---|---|---|
| | | | Precision | Recall | F1-score | Precision | Recall | F1-score |
| | EANN [42] | 0.827 | 0.847 | 0.812 | 0.829 | 0.807 | 0.843 | 0.825 |
| | MVAE [24] | 0.824 | 0.854 | 0.769 | 0.809 | 0.802 | 0.875 | 0.837 |
| | Spotfake [40] | 0.892 | 0.902 | **0.964** | **0.932** | 0.847 | 0.656 | 0.739 |
| | MVNN [45] | 0.846 | 0.809 | 0.857 | 0.832 | 0.879 | 0.837 | 0.858 |
| Weibo | SAFE [48] | 0.762 | 0.831 | 0.724 | 0.774 | 0.695 | 0.811 | 0.748 |
| | LIIMR [39] | 0.900 | 0.882 | 0.823 | 0.847 | 0.908 | **0.941** | **0.925** |
| | MCAN [44] | 0.899 | 0.913 | 0.889 | 0.901 | 0.884 | 0.909 | 0.897 |
| | CAFE [9] | 0.840 | 0.855 | 0.830 | 0.842 | 0.825 | 0.851 | 0.837 |
| | FND-CLIP | **0.907** | **0.914** | 0.901 | 0.908 | **0.914** | 0.901 | 0.907 |
| | RoBERTa-MWSS [37] | 0.820 | | - | - | 0.820 | - | - |
| | SAFE [48] | 0.874 | 0.851 | 0.830 | 0.840 | 0.889 | 0.903 | 0.896 |
| | Spotfake+ [38] | 0.846 | - | - | - | - | - | - |
| Politifact | TM [4] | 0.871 | - | - | - | 0.901 | - | - |
| | LSTM-ATT [27] | 0.832 | 0.828 | 0.832 | 0.830 | 0.836 | 0.832 | 0.829 |
| | DistilBert [1] | 0.741 | 0.875 | 0.636 | 0.737 | 0.647 | 0.880 | 0.746 |
| | CAFE [9] | 0.864 | 0.724 | 0.778 | 0.750 | 0.895 | 0.919 | 0.907 |
| | FND-CLIP | **0.942** | **0.897** | **0.897** | **0.897** | **0.960** | **0.960** | **0.960** |
| | RoBERTa-MWSS [37] | 0.800 | - | - | 0.800 | - | - | - |
| | SAFE [48] | 0.838 | 0.758 | 0.558 | 0.643 | 0.857 | 0.937 | 0.895 |
| | Spotfake+ [38] | 0.856 | - | - | - | - | - | - |
| Gossipcop | TM [4] | 0.842 | - | - | - | 0.896 | - | - |
| | LSTM-ATT [27] | 0.842 | **0.845** | **0.842** | **0.844** | 0.839 | 0.842 | 0.821 |
| | DistilBert [1] | 0.857 | 0.805 | 0.527 | 0.637 | 0.866 | **0.960** | 0.911 |
| | CAFE [9] | 0.867 | 0.732 | 0.490 | 0.587 | 0.887 | 0.957 | 0.921 |
| | FND-CLIP | **0.880** | 0.761 | 0.549 | 0.638 | **0.899** | 0.959 | **0.928** |

description were filtered out. If a news post contains a text with multiple associated images, we randomly select one image. Weibo is a widely used Chinese dataset in fake news detection. The training set contains 3, 749 real news and 3, 783 fake news, and the test set contains 1, 996 news. Politifact and Gossipcop datasets are two English datasets collected from the political and entertainment domains of FakeNewsNet [36] repository, respectively. Politifact contains 244 real news and 135 fake news in the training set and 75 real news and 29 news in the test set. Gossipcop contains 10, 010 training news, including 7, 974 real news and 2, 036 fake news. The test set contains 2, 285 real news and 545 fake news. Besides, while Twitter [6] is also a well-known multimodal dataset for FND, we find that it contains plenty of duplicated posts and over 10k posts host only 463 images. More importantly, more than 70% of tweets on Twitter dataset are related to a single event, which can easily lead to model overfitting. Therefore, we do not conduct experiments on Twitter.

**Baseline Methods.** For a fair and reproducible comparison, we have to be selective in choosing the baseline methods. First, we prefer methods that provide pre-trained models or source code publicly available. Second, the methods should follow a common evaluation protocol where the three datasets are used for training

and testing. Accordingly, we compare FND-CLIP with the following methods and provide a quick recap.

**EANN [42]**, which employs an auxiliary task of event classification to improve generalizability.

**MVAE [24]**, which uses a variational autoencoder to model representations between text and images for fake news detection.

**Spotfake [40]**, which uses VGG and BERT to respectively extract image and text features and concatenates them to classify.

**MVNN [45]**, which incorporates textual semantic features, visual tampering features, and similarity of textual and visual information in fake news detection.
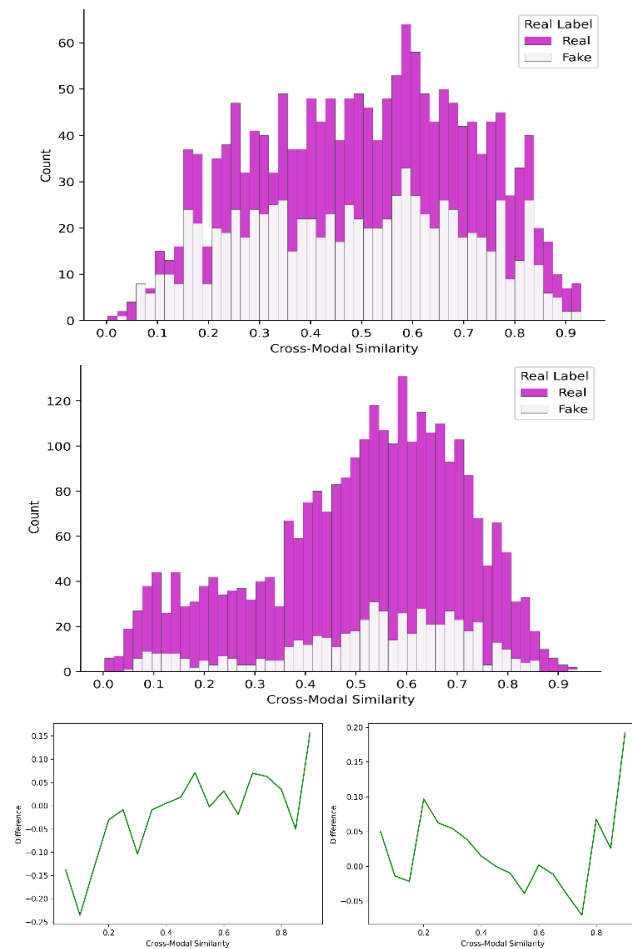
**SAFE [48]**, which fed the relevance between news textual and visual information into a classifier to detect fake news.

**LIIMR [39]**, which identifies and suppresses information from weaker modalities and extracts relevant information from the strong modality on a per-sample basis.

**MCAN [44]**, which stacks multiple co-attention layers to fuse the multimodal features.

**CAFE [9]**, which formulates an ambiguity-aware multimodal fake news detection method to adaptively aggregate unimodal features and cross-modal correlations.

**RoBERTa-MWSS [37]**, which exploits multiple weak signals from different sources from user and content engagements.

**Figure 3: Statistical analysis on cross-modal similarity of different news. The first and second row respectively show the counting of real/fake news according to cross-modal similarity on Weibo and Gossipcop. The third row shows the distance between the real news rate in each bin compared to the average rate on the corresponding dataset (left: Weibo. right: Gossipcop).**

**Spotfake+ [38]**, which is an improved version of Spotfake and can detect full length articles.

**TM [4]**, which utilizes lexical and semantic properties of both true and fake news text to detect fake news.

**LSTM-ATT [27]**, which builds a model based on XGBoost to detect full length fake news.

**DistilBert [1]**, which uses latent representations of news articles and user-generated content to guide model learning.

## 4.2 Performance Analysis

Table 1 shows the average precision, recall, and accuracy of FND-CLIP on three representative datasets. The results are promising, with over 90% average accuracy on Weibo and over 94% on Politifact, which indicates that the proposed method is a dependable

and robust fake news detection algorithm that can detect anomalies given multi-lingual and multi-domain news. Especially, the recall rates of real news on the three datasets are all above 0.9, and therefore FND-CLIP is less likely to classify a real news as fake.

To further conduct statistical analysis on how cross-modal similarity correlates with the attention score and how they vary given different news instances, Figure 3 shows the correlation between the CLIP-based cross-modal similarity score and the fake news ratio on Weibo and Gossipcop dataset. In row 1 and 2, we group all news in each dataset into several bins according to their similarity score, and find that a news is more likely to be real when the similarity score is high. In row 3, we calculate real news rates of each bin and subtract them with the average real news rate of the corresponding dataset. The curves show that real news rate goes up with the increasing ambiguity on Weibo, and first goes down then surges up on Gossipcop. Such statistical charateristics are useful for deep networks to identify fake news.

## 4.3 Comparison with State-of-the-arts

We further compare FND-CLIP with the above-mentioned state-of-the-art methods and the comparison results are presented in Table 1. '-' means the results are not available from the original paper. As shown in Table 1, FND-CLIP outperforms all the compared methods on the three datasets in terms of Accuracy, and achieves slightly lower than Spot on Weibo in Recall. FND-CLIP achieves the highest accuracy of 90.7%, 94.2%, and 88.0%, which surpasses 0.7%, 6.8%, and 1.3% over the state-of-the-art method, on the three real-world datasets, respectively. Besides, we rank either $1_{st}$ or $2_{nd}$ in precision, recall, and accuracy in all tests, which proves the effectiveness of FND-CLIP.

Many fake news detection methods, such as EANN and Spotfake, rely only on the fused features obtained by direct use of concatenating or attention mechanisms. However, these fused features cannot provide sufficient discrimination ability to classify fake news, because the text and image features separately extracted are not in the same semantic space and the correlation information of the text and image is not well-paid attention to during the fusion process. Therefore, the experimental results of these methods are unsatisfactory. CAFE uses cross-modal alignment to train encoders that can map texts and images into the same semantic space. By using the features fused from the alimented text and image features to classify, it achieves good experimental results, especially on the Politifact and Gossiopcop datasets. However, due to the limitation of the number of data sets and the rough label method for training labels, the encoding effect of the encoder is not optimal, and the semantic gap between text and image features is still significant. In addition, CAFE designs an ambiguity learning module to calculate a weight used for adaptively adjusting the calculation of different modalities. However, the weights for selecting unimodal or multimodal features are obtained by manual calculation, and cannot be further optimized by reverse gradient propagation, thus affecting the performance of the detection.

FND-CLIP outperforms most of the state-of-the-art methods, mainly due to the following reasons. First, the pre-trained CLIP encoders in FND-CLIP can generate semantically information-rich text and image features in the same semantic space, ensuring the

fused feature correctly reflects the correlation between text and image, and providing complementary information for the unimodal features. The modality-wise attention mechanism adaptively determines the weights of text, image, and fused features, avoiding the influence of invalid features on the representation ability of final features, and further improving the classification accuracy.

## 4.4 Ablation Studies

We explore the influence of the key components in FND-CLIP by evaluating the performance of the model with varied and partial setups. In each test, we remove different components and train the models from scratch. The compared variants of FND-CLIP are implemented as follows.

- FND-CLIP w/o A. We remove the modality-wise attention module and direct aggregate the three features to obtain final feature;
- FND-CLIP w/o F. We remove the fusion module and use two unimodal features to classify news;
- FND-CLIP w/o C. We remove all CLIP-related modules and only use BERT and ResNet to extract text and image features.
- FND-CLIP multimodal-only: We remove the unimodal feature extractor, BERT and ResNet, and only use CLIP fused feature as final feature;
- FND-CLIP image-only: We remove the all text-related features and only use image feature extracted by ResNet to classify;
- FND-CLIP text-only: We only use BERT-extracting feature to complete the detection task without any visual information.

**Effectiveness of Each Component.** First, we analyze the impact of different components in FND-CLIP for fake news detection. From the results shown in Table 2, we have the following observations: 1) FND-CLIP outperforms FND-CLIP w/o C, proving that CLIP can effectively provide discernable features for fake news detection task and significantly improve the accuracy of classification. Although only intra-modal features can be used for classification, the lack of interaction between modalities makes the final features lack the ability to represent the intrinsic relationship between images and texts. 2) FND-CLIP outperforms FND-CLIP w/o F, indicating that although the unimodal branches contain the CLIP-coded features, the fused feature reflecting the correlation of text and image provides effective multimodal information for classifier. Meanwhile, FND-CLIP w/o F outperforms FND-CLIP w/o C, indicating that the complement to unimodal features using CLIP-coded features is effective. 3) FND-CLIP outperforms FND-CLIP w/o A on Weibo and Gossipcop, indicating that modality-wise attention can help FND-CLIP adaptively weight useful modalities. FND-CLIP w/o A directly fuses the features of different modalities, which may cause the final feature be affected by invalid information from a modality. **Contributions from Different Modalities.** The second set of experiments is to evaluate the classification performance of different modalities in fake news detection. From Table 2, we draw some analysis as follows: 1) FND-CLIP image-only performs worst, especially on Gossipcop dataset, where the F1 score of fake news was almost zero, meaning that all news was judged real and the model had no classification ability at all. This shows that in fake news detection, simple visual information provides fewer classification

clues than other modalities. 2) FND-CLIP multimodal-only achieves accuracy of 81.7%, 90.3%, and 86.2% on Weibo, Politifact, and Gossipcop datasets respectively, but performs worse than FND-CLIP text-only on Weibo and Gossipcop datasets, indicating that the correlation information of images and texts can be used to classify fake news. However, the classification ability of fused feature is limited because news itself has modal irrelevance and ambiguity. In addition, CLIP-based fused features focus on the semantics of the text, while the BERT-based text features also extract emotional features that are helpful for fake news detection. 3) FND-CLIP text-only achieves the second-best results, indicating that only using text feature can basically complete the classification task for fake news. However, FND-CLIP outperforms FND-CLIP text-only, proving the visual feature can supplement classification information and the correct use of multimodal features is superior to using only unimodal features in fake news detection.

## 4.5 T-SNE Visualizations

In Figure 4, we further analyze the proposed method using t-SNE [41] visualizations of the features before classifier that are learned by FND-CLIP, CAFE, and also the proposed method with partial settings such as FND-CLIP w/o C, FND-CLIP w/o A, FND-CLIP text-only, and FND-CLIP image-only on the test dataset of Weibo in Figure 4.

The dots with the same color mean that they are within the same label. From Figure 4 we can see that the boundary of different label dots in FND-CLIP is more pronounced than that in CAFE, FND-CLIP w/o C, and FND-CLIP w/o A, revealing that the extracted features in FND-CLIP are more discriminative than those in CAFE and the CLIP-related modules and modality-wise attention are useful for improving the classification ability of FND-CLIP.

In addition, by comparing Figure 4a, Figure 4d, Figure 4e, and Figure 4f, we can see that image features alone are not enough for classification, which indicates that the image itself does not have classification ability. The effect of text-only is much better than that of image-only, Proving that the text features play a leading role in fake news detection, but there are still many sample dots that cannot be distinguished. FND-CLIP w/o C, which contains both text and image features, has a more obvious boundary of the dots than FND-CLIP text-only, indicating that different modalities have complementary information. In addition, the separation degree of the sample dots in Figure 4a is higher than that in Figure 4d, indicating that the multimodal features based on CLIP can improve the representation ability of the final features.

## 5 CONCLUSIONS

In this paper, we present a novel multimodal fake news detection method called FND-CLIP, which uses CLIP to extract aligned multimodal features and guide the learning of network for different modalities. In addition, we introduce modality-wise attention to adaptively determine the weights of text, image, and fused features. It can avoid introducing noisy and redundant features during feature fusion, which further improve the classification accuracy. We conduct comprehensive experiments on several well-known FND datasets. The results show that using CLIP for multimodal feature generation can well collaborate with the unimodal features

**Table 2: Ablation study on the architecture design and different features of FND-CLIP on three datasets. The entire FND-CLIP achieved the highest accuracy and F1-score, demonstrating that every module in the architecture of our method is effectiveness and every modality is effectively utilized.**
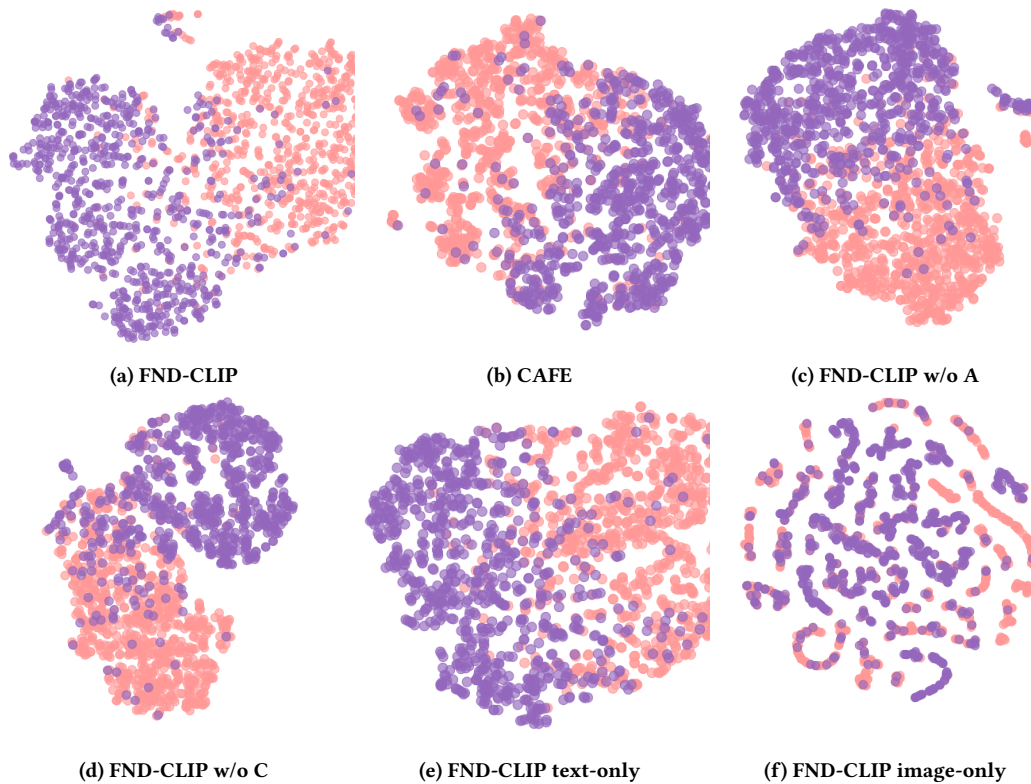
| | Method | Accuracy | Fake News | | | Real News | | |
|---|---|---|---|---|---|---|---|---|
| | | | Precision | Recall | F1-score | Precision | Recall | F1-score |
| Weibo | FND-CLIP multimodal-only | 0.817 | 0.899 | 0.718 | 0.798 | 0.761 | 0.917 | 0.832 |
| | FND-CLIP image-only | 0.796 | 0.862 | 0.711 | 0.779 | 0.750 | 0.884 | 0.811 |
| | FND-CLIP text-only | 0.872 | 0.906 | 0.833 | 0.868 | 0.842 | 0.911 | 0.875 |
| | FND-CLIP w/o C | 0.874 | 0.895 | 0.851 | 0.872 | 0.855 | 0.898 | 0.876 |
| | FND-CLIP w/o F | 0.893 | 0.925 | 0.857 | 0.890 | 0.864 | 0.929 | 0.895 |
| | FND-CLIP w/o A | 0.897 | **0.936** | 0.855 | 0.893 | 0.863 | **0.940** | 0.900 |
| | FND-CLIP | **0.907** | 0.914 | **0.901** | **0.908** | **0.901** | 0.914 | **0.907** |
| Politifact | FND-CLIP multimodal-only | 0.903 | 0.807 | 0.862 | 0.833 | 0.944 | 0.919 | 0.932 |
| | FND-CLIP image-only | 0.748 | 0.600 | 0.310 | 0.409 | 0.773 | 0.919 | 0.840 |
| | FND-CLIP text-only | 0.903 | 0.913 | 0.724 | 0.808 | 0.900 | **0.973** | 0.935 |
| | FND-CLIP w/o C | 0.893 | 0.875 | 0.724 | 0.793 | 0.899 | 0.960 | 0.928 |
| | FND-CLIP w/o F | 0.903 | 0.880 | 0.759 | 0.815 | 0.910 | 0.960 | 0.934 |
| | FND-CLIP w/o A | **0.942** | **0.926** | 0.862 | 0.893 | 0.947 | **0.973** | **0.960** |
| | FND-CLIP | **0.942** | 0.897 | **0.897** | **0.897** | **0.960** | 0.960 | **0.960** |
| Gossipcop | FND-CLIP multimodal-only | 0.862 | 0.708 | 0.484 | 0.575 | 0.886 | 0.952 | 0.918 |
| | FND-CLIP image-only | 0.814 | **1.000** | 0.033 | 0.064 | 0.813 | **1.000** | 0.897 |
| | FND-CLIP text-only | 0.871 | 0.741 | 0.508 | 0.603 | 0.891 | 0.958 | 0.923 |
| | FND-CLIP w/o C | 0.870 | 0.745 | 0.494 | 0.594 | 0.888 | 0.960 | 0.923 |
| | FND-CLIP w/o F | 0.874 | 0.723 | 0.562 | 0.632 | 0.901 | 0.949 | 0.924 |
| | FND-CLIP w/o A | 0.873 | 0.715 | **0.567** | 0.633 | **0.902** | 0.946 | 0.923 |
| | FND-CLIP | **0.880** | 0.761 | 0.549 | **0.638** | 0.899 | 0.959 | **0.928** |

extracted by ResNet and BERT in mining crucial features for fake news detection. More importantly, FND-CLIP outperforms many of the state-of-the-art methods in multimodal fake news detection.

Aside from the performance gain of FND-CLIP, the outputs are still in the form of binary value that predicts either "real" or "fake", which cannot somehow explain why the news is predicted fake and which elements in the news are most suspicious and abnormal. In future works, we head towards developing more explainable fake news detection systems that can provide reasons why a given news is predicted as real or fake.

## REFERENCES

[1] Liesbeth Allein, Marie-Francine Moens, and Domenico Perrotta. 2021. Like Article, Like Audience: Enforcing Multimodal Correlations for Disinformation Detection. *arXiv preprint arXiv:2108.13892* (2021).

[2] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. 2018. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence* 41, 2 (2018), 423–443.

[3] Gaurav Bhatt, Aman Sharma, Shivam Sharma, Ankush Nagpal, Balasubramanian Raman, and Ankush Mittal. 2018. Combining neural, statistical and external features for fake news stance identification. In *Companion Proceedings of the The Web Conference 2018.* 1353–1357.

[4] Bimal Bhattarai, Ole-Christoffer Granmo, and Lei Jiao. 2021. Explainable Tsetlin Machine framework for fake news detection with credibility score assessment. *arXiv preprint arXiv:2105.09114* (2021).

[5] Tian Bian, Xi Xiao, Tingyang Xu, Peilin Zhao, Wenbing Huang, Yu Rong, and Junzhou Huang. 2020. Rumor detection on social media with bi-directional graph convolutional networks. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 549–556.

[6] Christina Boididou, Symeon Papadopoulos, Markos Zampoglou, Lazaros Apostolidis, Olga Papadopoulou, and Yiannis Kompatsiaris. 2018. Detection and visualization of misleading content on Twitter. *International Journal of Multimedia Information Retrieval* 7, 1 (2018), 71–86.

[7] 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* (2020), 141–161.

[8] Xinru Chen, Chengbo Dong, Jiaqi Ji, Juan Cao, and Xirong Li. 2021. Image Manipulation Detection by Multi-View Multi-Scale Supervision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision.* 14185–14193.

[9] Yixuan Chen, Dongsheng Li, Peng Zhang, Jie Sui, Qin Lv, Lu Tun, and Li Shang. 2022. Cross-modal Ambiguity Learning for Multimodal Fake News Detection. In *Proceedings of the ACM Web Conference 2022.* 2897–2905.

[10] Marcos V Conde and Kerem Turgutlu. 2021. CLIP-Art: contrastive pre-training for fine-grained art classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 3956–3960.

[11] Nadia K Conroy, Victoria L Rubin, and Yimin Chen. 2015. Automatic deception detection: Methods for finding fake news. *Proceedings of the association for information science and technology* 52, 1 (2015), 1–4.

[12] Corentin Dancette, Remi Cadene, Damien Teney, and Matthieu Cord. 2021. Beyond question-based biases: Assessing multimodal shortcut learning in visual question answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision.* 1574–1583.

[13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[14] Han Guo, Juan Cao, Yazi Zhang, Junbo Guo, and Jintao Li. 2018. Rumor detection with hierarchical social attention network. In *Proceedings of the 27th ACM international conference on information and knowledge management.* 943–951.

[15] Bing Han, Xiaoguang Han, Hua Zhang, Jingzhi Li, and Xiaochun Cao. 2021. Fighting fake news: two stream network for deepfake detection via learnable SRM. *IEEE Transactions on Biometrics, Behavior, and Identity Science* 3, 3 (2021), 320–331.

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 770–778.

[17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 770–778.

**Figure 4: T-SNE visualizations of the features before classifier that are learned by FND-CLIP , CAFE, FND-CLIP w/o C, FND-CLIP w/o A, FND-CLIP text-only, and FND-CLIP image-only on the test dataset of Weibo.**

[18] Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415* (2016).

[19] Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7132–7141.

[20] Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*. 448–456.

[21] Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. What does BERT learn about the structure of language?. In *ACL 2019-57th Annual Meeting of the Association for Computational Linguistics*.

[22] Zhiwei Jin, Juan Cao, Han Guo, Yongdong Zhang, and Jiebo Luo. 2017. Multimodal fusion with recurrent neural networks for rumor detection on microblogs. In *Proceedings of the 25th ACM international conference on Multimedia*. 795–816.

[23] Gregory Johnson. 2012. Google Translate http://translate. google. com. *Technical Services Quarterly* 29, 2 (2012), 165–165.

[24] Dhruv Khattar, Jaipal Singh Goud, Manish Gupta, and Vasudeva Varma. 2019. Mvae: Multimodal variational autoencoder for fake news detection. In *The world wide web conference*. 2915–2921.

[25] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[26] Peiguang Li, Xian Sun, Hongfeng Yu, Yu Tian, Fanglong Yao, and Guangluan Xu. 2021. Entity-Oriented Multi-Modal Alignment and Fusion Network for Fake News Detection. *IEEE Transactions on Multimedia* (2021).

[27] Jun Lin, Glenna Tremblay-Taylor, Guanyi Mou, Di You, and Kyumin Lee. 2019. Detecting fake news articles. In *2019 IEEE International Conference on Big Data (Big Data)*. IEEE, 3021–3025.

[28] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems* 32 (2019).

[29] Qiong Nan, Juan Cao, Yongchun Zhu, Yanyan Wang, and Jintao Li. 2021. MD-FEND: Multi-domain Fake News Detection. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 3343–3347.

[30] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2021. Glide: Towards photorealistic

[31] Martin Potthast, Johannes Kiesel, Kevin Reinartz, Janek Bevendorff, and Benno Stein. 2017. A stylometric inquiry into hyperpartisan and fake news. *arXiv preprint arXiv:1702.05638* (2017).

[32] 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 *Proceedings of the 29th ACM International Conference on Multimedia*. 1212–1220.

[33] 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)*. IEEE, 518–527.

[34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*. PMLR, 8748–8763.

[35] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683* (2019).

[36] Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. 2020. Fakenewsnet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media. *Big data* 8, 3 (2020), 171–188.

[37] Kai Shu, Guoqing Zheng, Yichuan Li, Subhabrata Mukherjee, Ahmed Hassan Awadallah, Scott Ruston, and Huan Liu. 2020. Leveraging multi-source weak social supervision for early detection of fake news. *arXiv preprint arXiv:2004.01732* (2020).

[38] Shivangi Singhal, Anubha Kabra, Mohit Sharma, Rajiv Ratn Shah, Tanmoy Chakraborty, and Ponnurangam Kumaraguru. 2020. Spotfake+: A multimodal framework for fake news detection via transfer learning (student abstract). In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 13915–13916.

image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741* (2021).

[39] Shivangi Singhal, Tanisha Pandey, Saksham Mrig, Rajiv Ratn Shah, and Ponnurangam Kumaraguru. 2022. Leveraging Intra and Inter Modality Relationship for Multimodal Fake News Detection. (2022).

[40] Shivangi Singhal, Rajiv Ratn Shah, Tanmoy Chakraborty, Ponnurangam Kumaraguru, and Shin'ichi Satoh. 2019. Spotfake: A multi-modal framework for fake news detection. In *2019 IEEE fifth international conference on multimedia big data (BigMM)*. IEEE, 39–47.

[41] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).

[42] 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*. 849–857.

[43] Tianyi Wei, Dongdong Chen, Wenbo Zhou, Jing Liao, Zhentao Tan, Lu Yuan, Weiming Zhang, and Nenghai Yu. 2021. Hairclip: Design your hair by text and reference image. *arXiv preprint arXiv:2112.05142* (2021).

[44] Yang Wu, Pengwei Zhan, Yunjian Zhang, Liming Wang, and Zhen Xu. 2021. Multimodal Fusion with Co-Attention Networks for Fake News Detection. In

[45] Junxiao Xue, Yabo Wang, Yichen Tian, Yafei Li, Lei Shi, and Lin Wei. 2021. Detecting fake news by exploring the consistency of multimodal data. *Information Processing & Management* 58, 5 (2021), 102610.

[46] Keren Ye and Adriana Kovashka. 2021. A case study of the shortcut effects in visual commonsense reasoning. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 35. 3181–3189.

[47] Xueyao Zhang, Juan Cao, Xirong Li, Qiang Sheng, Lei Zhong, and Kai Shu. 2021. Mining dual emotion for fake news detection. In *Proceedings of the Web Conference 2021*. 3465–3476.

[48] Xinyi Zhou, Jindi Wu, and Reza Zafarani. 2020. SAFE: Similarity-Aware Multimodal Fake News Detection. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 354–367.

[49] Arkaitz Zubiaga, Ahmet Aker, Kalina Bontcheva, Maria Liakata, and Rob Procter. 2018. Detection and resolution of rumours in social media: A survey. *ACM Computing Surveys (CSUR)* 51, 2 (2018), 1–36.

*Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021.* 2560–2569.