

# Deep Visual–Linguistic Fusion Network Considering Cross-Modal Inconsistency for Rumor Detection

Yang Yang<sup>1\*</sup>, Ran Bao<sup>2</sup>, Weili Guo<sup>1</sup>, De-Chuan Zhan<sup>2</sup>, Yilong Yin<sup>3</sup> & Jian Yang<sup>1</sup>

<sup>1</sup>Nanjing University of Science and Technology, Nanjing 210094, China  
{yyang,wlguo,csjyang}@njut.edu.cn;

<sup>2</sup>National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China  
baoranx@gmail.com, zhandc@nju.edu.cn;

<sup>3</sup>School of Software, Shandong University, Shandong 250000, China  
ylyin@sdu.edu.cn

## Abstract

With the development of the Internet, users can freely publish posts on various social media platforms, which offers great convenience for keeping abreast of the world. However, posts usually carry many rumors, which require plenty of manpower for monitoring. Owing to the success of modern machine learning techniques, especially deep learning models, we tried to detect rumors as a classification problem automatically. Early attempts have always focused on building classifiers relying on image or text information, i.e., single modality in posts. Thereafter, several multimodal detection approaches employ an early or late fusion operator for aggregating multiple source information. Nevertheless, they only take advantage of multimodal embeddings for fusion and ignore another important detection factor, i.e., the intermodal inconsistency between modalities. To solve this problem, we develop a novel deep visual–linguistic fusion network (DVLFN) considering cross-modal inconsistency, which detects rumors by comprehensively considering modal aggregation and contrast information. Specifically, the DVLFN first utilizes visual and textual deep encoders, i.e., Faster R-CNN and Bidirectional Encoder Representations from Transformers, to extract global and regional embeddings for image and text modalities. Then, it predicts posts' authenticity from two aspects: 1) intermodal inconsistency, which employs the Wasserstein distance to efficiently measure the similarity between regional embeddings of different modalities, and 2) modal aggregation, which experimentally employs the early fusion to aggregate two modal embeddings for prediction. Consequently, the DVLFN can compose the final prediction based on the modal fusion and inconsistency measure. Experiments are conducted on three real-world multimedia rumor detection datasets collected from Reddit, GoodNews, and Weibo. The results validate the superior performance of the proposed DVLFN.

**Keywords** Multimodal Learning, Wasserstein Distance, Rumor Detection

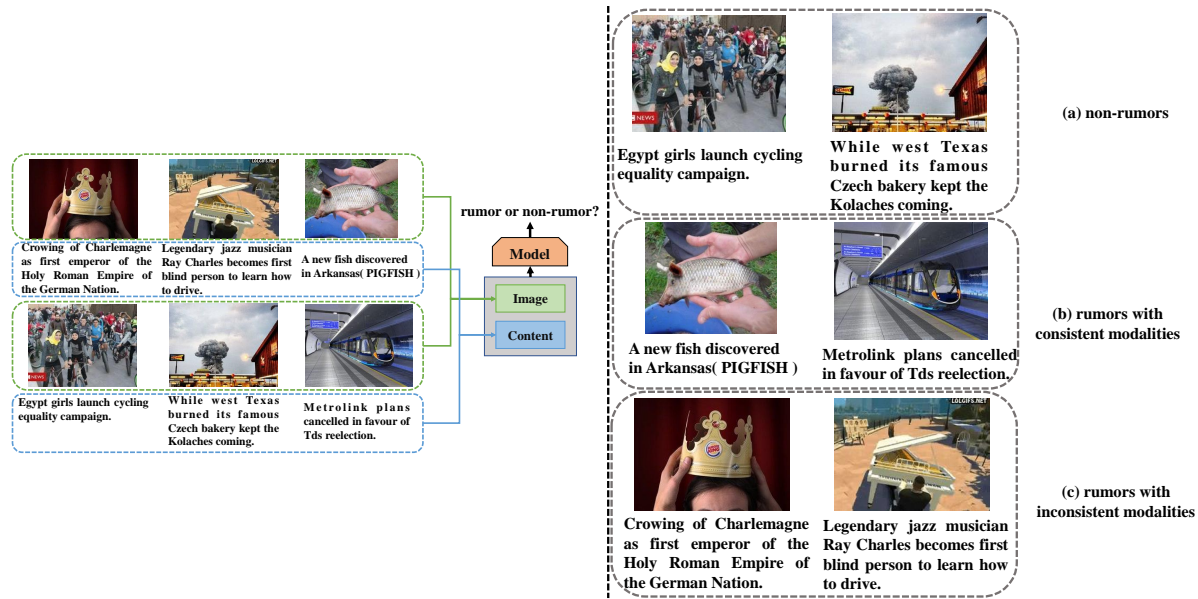
**Citation** Y Yang, R Bao, W.L Guo, D.C Zhan, Y.L Yin, J Yang.  
. Sci China Inf Sci, for review

## 1 Introduction

With the development of social media platforms, such as Twitter, Instagram, and Weibo, users can easily read and publish real-time posts. This advancement brings great convenience for governments and organizations to release real-time news and for users to keep abreast of surrounding events. For example, in the early stage of the COVID-19 pandemic, several reports on the global epidemic progress and related prevention suggestions were widely spread and rapidly forwarded on Weibo, Twitter, and other platforms. This greatly facilitated the control and prevention of the epidemic and had great social significance. However, various rumors emerged due to the convenience and openness of social media and caused a serious negative social impact [1,2]. For example, rumors about presidential candidates emerged in the 2016 U.S. presidential election, which garnered significant attention in media and policy

1) Yang Yang and Jian Yang are with PCA Lab, Key Lab of Intelligent Perception and Systems for High-Dimensional Information of Ministry of Education, and Jiangsu Key Lab of Image and Video Understanding for Social Security, School of Computer Science and Engineering, Nanjing University of Science and Technology. Yang Yang is also with the Key Laboratory of Computer Network and Information Integration (Southeast University), Ministry of Education. State Key Lab. for Novel Software Technology, Nanjing University, P.R. China.

\* Corresponding author (email: yyang@njut.edu.cn)



**Figure 1** (Best viewed in color when zoomed in.) Multimodal detection framework (left) and intuition of the DVLFN (right). Multimodal detection framework always leverages possible visual–linguistic fusion for a better performance. Meanwhile, rumors can be detected from two points: (b) Rumors with consistent modalities. The embeddings of the image and content are consistent, and we can verify the authenticity from visual or textual information. (c) Rumors with inconsistent modalities. Several rumors are difficult to verify the authenticity from any modality but are much easier to verify from the inconsistency degree of an image–text pair.

circles. Some journalists even claimed that the results of the 2016 election were a consequence of the rumors spread. In addition, rumors that amoxicillin can treat COVID-19 received a lot of attention and were widely reposted within a short time, which caused item shortage and public panic. Numerous fake posts on social media platforms may also greatly mislead consumers and even directly affect financial activities. These rumors have raised fears by misinforming users and corroding our society [3]. To solve this problem, fact-checking websites have been built to clarify rumors, such as snopes.com<sup>1)</sup> and politifact.com<sup>2)</sup>. However, they always rely on manual selection to detect rumors, which has obvious limitations on efficiency considering the large number of posts on social media. Therefore, the automatic detection of rumors is crucial to increasing the credibility of real-time posts.

As a matter of fact, posts on modern social media usually contain text and image information. With the rapid progression in computer vision [4,5] and natural language processing [6–8], we have confidence in detecting rumors from mass posts using artificial intelligence techniques. Initial approaches [9–12] usually concentrated on mining the text modality. They mainly employed handcrafted semantic features, such as term frequency–inverse document frequency, bag-of-words (BOW), or statistical features, to verify rumors, but these shallow features have a limited prediction ability. Inspired by the great success of deep neural networks, deep models have been recently exploited and demonstrated state-of-the-art performance. For instance, [13] employed recurrent neural networks (RNNs), i.e., long short-term memory (LSTM) and gated recurrent unit (GRU), to learn hidden features from text content for detecting rumors. [14] adopted text convolutional neural networks (CNNs) for obtaining local and global features from relevant posts. Meanwhile, image modality can depict an event. Hence, researchers have tried several deep visual approaches to analyze image modality, i.e., whether they are artificially generated or not. For instance, [15] extracted forensic features to evaluate the authority of the attached images. [16] fused multiple domain visual information for detection. Furthermore, inspired by the generative adversarial networks (GANs) [17], several approaches [18–20] modeled rumor detection as a text-based or image-based GAN-style framework to generate confusing training examples for enhancing the representation learning of rumor-indicative patterns. However, these methods focus on a single modality, i.e., image or text modality, but they ignore the fact that an effective multimodal fusion can obtain more discriminative embeddings and better predictions.

Based on this idea, researchers have attempted to combine two modal information in detection using multimodal fusion techniques. As shown in the left side of Figure 1, [21–24] deployed relative deep multimodal fusion networks for the rumor detection task. They usually directly fused multimodal deep embeddings or learned consistent embeddings and adopted a corresponding structural regularization or multitask loss for joint optimization. The basic

1) <https://www.snopes.com/>

2) <https://www.politifact.com/>

assumption behind them is that either images or texts can verify rumors. These multimodal fusion methods mainly integrated two modalities' information, leaving the intermodal information without consideration. Figure 2 shows the inconsistency distribution (calculated with Wasserstein matching) on Fakeddit dataset [25] examples. We find that most non-rumors have small inconsistency values, whereas rumors have large inconsistency values. Therefore, we can actually detect rumors from two aspects: 1) Rumors with a consistent image–text pair, i.e., information coming from an image and text is consistent, and either an image or text can reflect the falsity. For example, as shown in Figure 1 (b), the image and text depict the same event, whereas we can detect a rumor from the image/text or both. 2) Rumors with an inconsistent image–text pair, i.e., it is difficult to distinguish rumors from any modal information, but we can confirm rumors by detecting the mismatching of two modal contents. For example, as indicated in Figure 1 (c), the text descriptions and attached images are mismatching. To sum up, existing multimodal methods aim to learn more discriminative embeddings by fusing all modal information, thus detecting rumors. Such an operator is applicable for the first type of rumors, but it cannot effectively detect the second type of rumors.

To this end, we propose a novel deep visual–linguistic fusion network (DVLFN), which includes intra-modal and intermodal evaluations, to comprehensively fuse multimodal information. Specifically, the DVLFN employs independently visual- and textual-based deep encoders to learn regional and global embeddings and then develops two evaluation criteria: 1) Modal aggregation. Through experimental verification, we use the early fusion (i.e., embedding fusion) to concatenate each modal embedding for prediction. 2) Intermodal inconsistency. The DVLFN calculates the Wasserstein distance between two modal region embeddings as an intermodal inconsistency. Consequently, the DVLFN combines the two measurements for the final prediction. Extensive experiments on three real-world datasets validate the superiority of using the DVLFN. Particularly, in terms of the F1 score, we achieve an absolute 2.7% improvement over the baselines on the NeuralNews dataset [26]. In summary, the contributions of this paper are as follows:

- We propose a novel DVLFN, which considers the extra intermodal inconsistency for rumor detection.
- We utilize visual and textual deep encoders to acquire global and regional embeddings and employ extra Wasserstein matching to comprehensively predict authenticity.
- In the experiments, our approach improves the performance on three real-world datasets, which validate that modal aggregation and intermodal knowledge are effective for enhancing the detector.

**Discussion.** Existing multimodal detection methods usually focus on designing effective fusion approaches for a good detection, and the method considering modal inconsistency only adopts manually designed inconsistent features or coarse-grained inconsistency. Actually, modal inconsistency is an effective criterion for rumor detection. To qualitatively detect modal inconsistency, we introduce the fine-grained inconsistent measure by Wasserstein matching, which aims to find subtler differences rather than simple aggregations. As a result, with modal fusion and modal inconsistency, we can acquire a good detection performance.

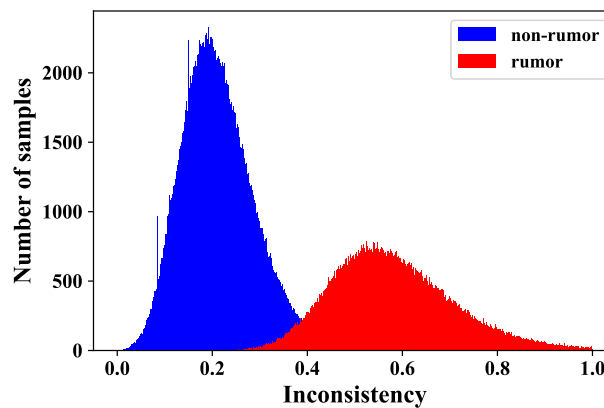


Figure 2 Inconsistency distribution of the Fakeddit dataset.

## 2 Related Works

This paper aims to detect rumors by comprehensively considering the modal fusion and inconsistency information of posts. Therefore, our work is related to single-modal detection, multimodal detection, and cross-modal learning.

## 2.1 Single-Modal Detection

As a text is the main component of posts, traditional methods always utilize extracted semantic representations based on texts for rumor detection. For example, [9, 27] adopted BOW features from texts to extract relationships among posts and detect rumors. [12] employed a latent Dirichlet allocation (LDA) model to represent abstract semantic features for detection. However, these methods are limited to manually crafted features, which affect the detection performance. With the development of deep neural networks, several research works have turned to utilizing end-to-end text deep networks for the detection task. For example, [13] utilized the RNN model to learn high-level semantic embeddings of posts and predicted the possibility of authenticity. In addition to textual information, statistical features from texts, such as count of words, punctuation, hashtag topics(#), mentions(@), and URLs, were also utilized as auxiliary features [11, 28], which can capture the prominent statistical information for assisting detection. For example, [29] combined statistical features for detection, and [30] incorporated social contexts as auxiliary information into a hierarchical neural network via the attention mechanism for detection. Moreover, considering complementary visual information from attached images, several studies focused on distinguishing rumors by measuring the image quality. For example, [15] proposed to use visual forensic features for detection, and [31] designed a CNN-based network to automatically capture the complex patterns of forged images in the frequency domain. Recently, inspired by the GAN [17], several works combined text- or image-based GANs to learn more robust embeddings. For example, [18] extracted co-occurrence matrices on three-color channels in the pixel domain and trained a robust model using the deep CNN framework. [19] proposed a text-based GAN-style approach, where the generator produced conflicting posts to pressurize the discriminator and learn strong rumor-indicative representations. However, all these methods are based on a single modality and cannot effectively fuse multimodal information contained in posts.

## 2.2 Multimodal Detection

Recent studies aim to verify the credibility of multimedia posts by fusing multimodal information, i.e., texts and images. For example, [32] proposed the verifying multimedia task as a part of the MediaEval benchmark in 2015 and 2016. Multimodal deep neural networks are also proposed to fuse multimodal embeddings, in which CNNs or RNNs are always employed as a basic model for image/text modality. For example, [24] first incorporated deep neural networks for rumor detection by concatenating deep multimodal embeddings. [23] proposed an end-to-end event adversarial neural network to detect emerged rumors based on learning invariant multimodal embeddings. [22] trained a multimodal variational autoencoder jointly with a rumor detector to learn shared embeddings for texts and images. These methods always directly fuse multimodal embeddings or learn potentially consistent feature embeddings for detection, but they may receive interferences from noise information of indecipherable modality in inconsistent cases.

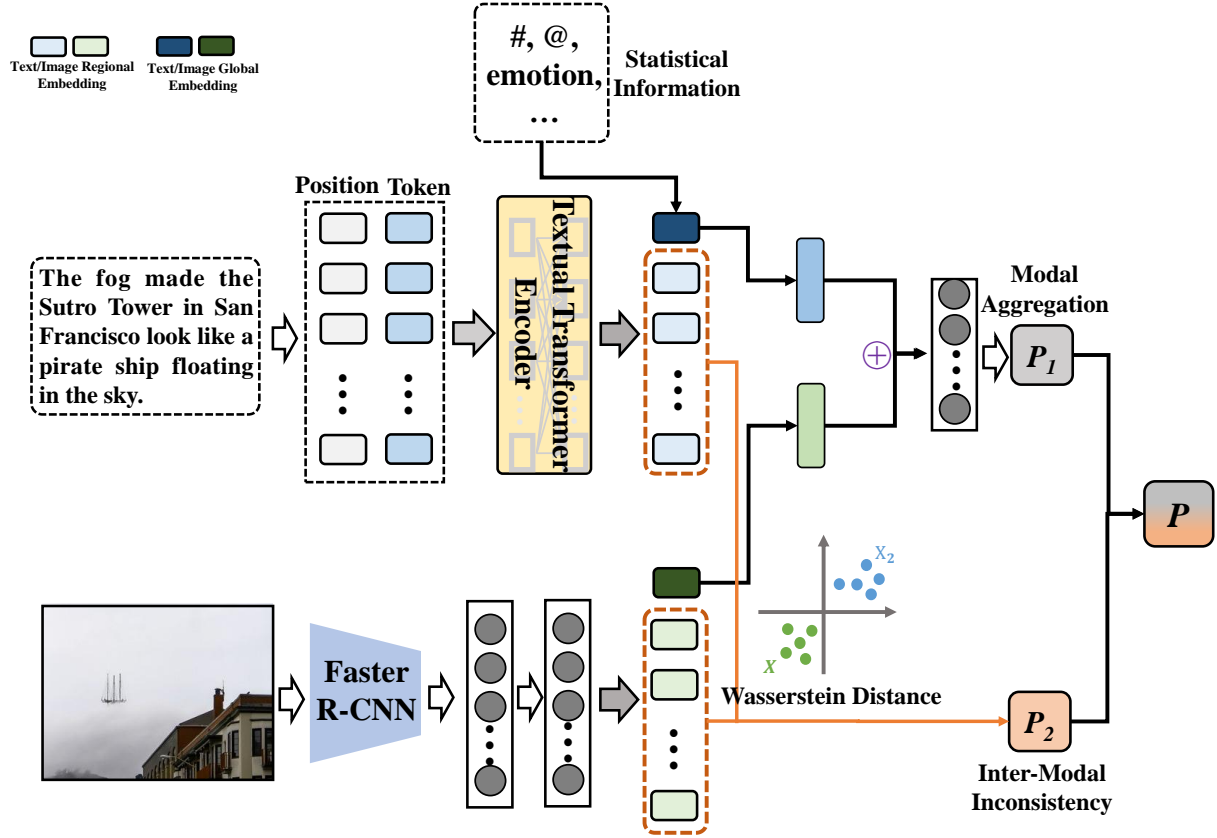
## 2.3 Cross-Modal Learning

Our work employs the inconsistency between two modalities, which is related to cross-modal learning. Currently, cross-modal learning aims to bridge connections among different modalities and has many applications, such as image captioning [33], cross-modal retrieval [34, 35], and visual question answering [36]. These tasks mainly utilize cross-modal consistency to construct the cross-modal generator or learn consistent embeddings. For example, state-of-the-art approaches in bidirectional image–sentence retrieval [37, 38] have leveraged visual–linguistic consistency to learn consistent embeddings and achieved great success on standard datasets, such as MSCOCO [39] and Flickr30K [40]. These methods usually adopted the modern distance or similarity measure to calculate the consistency degree between two modal global or regional embeddings. In this paper, we aim to utilize the cross-modal inconsistency in reverse. [26, 41, 42] mentioned the cross-modal inconsistency for detection, but these methods focused on manually designed inconsistent features or global inconsistency measure.

The remainder of this paper is organized as follows: Section 3 presents the proposed method, including the model, solution, and extension. Section 4 shows the experimental results on three rumor detection datasets under different settings. Section 5 concludes this paper.

## 3 Proposed Method

We describe the details of our proposed DVLFN method in this section. The main goal of the DVLFN is to build a deep visual–linguistic fusion model for rumor detection, which not only integrates the intermodal inconsistency



**Figure 3** (Best viewed in color when zoomed in.) Illustration of the proposed DVLFN. It has three parts: 1) Textual model, i.e., textual transformer, fuses textual and statistical information to acquire semantic global and regional embeddings. 2) Visual model, i.e., Faster R-CNN, learns visual global and regional embeddings. 3) Prediction module includes intermodal inconsistency by calculating the Wasserstein distance and multimodal fusion through a generalized cross-entropy.

measure but also considers different modal fusions.

### 3.1 Notation

We first describe notations used throughout this paper. Suppose there exist  $N$  labeled posts with multimodal information, i.e.,  $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$ ,  $\mathbf{x}_i = \{\mathbf{x}_i^v, \mathbf{x}_i^w\}$  has two modalities, where superscript  $v$  denotes image modality and  $w$  denotes text modality. Note that  $\mathbf{x}_i^w$  includes statistical information  $\mathbf{x}_i^s$ .  $\mathbf{y}_i \in \{0, 1\}^C$ , where  $C = 2$ ,  $y_i = 1$  denotes a rumor; otherwise, it is a non-rumor.

### 3.2 Model Pipeline

Without any loss of generality, the DVLFN can be divided into two important modules: 1) Embedding encoders, which learn discriminative embeddings for each modality. The key challenge is the selection of a multimodal encoder. In other words, we need to consider building an independent encoder for each modality or an interactive shared encoder for two modalities. In this study, we choose independent encoders for the following reasons: a) Influence of embedding learning. The importance of different modalities is different, and thus a shared model may lead to a weak modality (i.e., modalities with weak classification performance), negatively influencing the embedding learning of the strong modality (i.e., modalities with strong classification performance) [43]. b) Influence of the inconsistency measure. We focus on the inconsistency prediction, whereas the shared model actually tends to learn consistent embeddings, which may have a negative impact. Therefore, we develop visual and textual deep encoders, and we compare the two methods in the ablation study. 2) Classifier. We directly fuse global embeddings and then acquire the prediction with the fully connected (FC) network. Meanwhile, we adopt Wasserstein matching to measure the inconsistency calculation by region embeddings. The final prediction is constituted by the two results. The whole framework is shown in Figure 3. The image extracts raw global and regional features through the Faster R-CNN, whereas the text adopts the textual transformer encoder to learn global and regional embeddings. Ultimately, we aggregate the fused modal prediction and intermodal inconsistency for the final detection.

### 3.3 Basic Networks

#### 3.3.1 Textual Encoder

In the upper-left corner of Figure 3, we first present the textual encoder  $g_w$  for extracting text semantic embedding. Considering the success of BERT [44], we employ the transformer encoder for  $g_w$ . The input text is first tokenized into a token sequence according to WordPieces [45], followed by the standard BERT preprocessing method. We also add the special token  $[CLS]$  for learning global embeddings. Accordingly, the encoder can be represented as

$$\begin{aligned}\hat{\mathbf{e}}^w &= BERT(\bar{\mathbf{e}}^w) \\ \bar{\mathbf{e}}^w &= LN(\psi_{we}(\mathbf{x}^w) + \psi_{wl}(\mathbf{p}^w)),\end{aligned}\quad (1)$$

where  $\hat{\mathbf{e}}^w$  denotes output embeddings and  $\bar{\mathbf{e}}^w$  represents the input of the BERT encoder. The representation for each sub-word token  $\bar{\mathbf{e}}^w$  is obtained by summing up its word embedding and position embedding, followed by an LN. In detail, the instance can be represented as  $\mathbf{x}^w = \{\mathbf{x}_l^w\}_{l=1}^{L_w}$ , where  $L_w$  denotes the instance length. The raw text representations  $\mathbf{x}^w$  (initialized from input tokens) and position features  $\mathbf{p}^w$  are fed through FC networks, i.e.,  $\psi_{we}, \psi_{wl}$ , which project them into the same embedding space. Then, we sum the two features and use the LN to obtain  $\bar{\mathbf{e}}^w$ . We adopt BERT as the encoder to transform the embedding features layer by layer with adaptive attention weights.

The BERT encoder is a stack of identical blocks, which consists of multi-head attention and feed-forward layers, both wrapped in residual adds [8]. In detail, the  $k$ -th layer output embeddings can be represented as  $\mathbf{z}^k = \{\mathbf{z}_l^k\}_{l=1}^{L_w}$ , where  $L_w$  denotes the input sequence length. We use sequence mask to uniform the length of tokens to  $L_w$ .  $\mathbf{z}^0$  is set as the input  $\bar{\mathbf{e}}^w$ . Then, the feature of the element  $l$  in the  $(k+1)$ -th layer  $\mathbf{z}_l^{k+1}$  can be computed by

$$\begin{aligned}\bar{\mathbf{h}}_l^{k+1} &= \sum_{m=1}^M W_m^{k+1} \left( \sum_{j=1}^L A_{l,j}^m \cdot V_m^{k+1} \mathbf{z}_j^k \right), \\ \mathbf{h}_l^{k+1} &= LN(\mathbf{z}_l^k + \bar{\mathbf{h}}_l^{k+1}), \\ \bar{\mathbf{z}}_l^{k+1} &= W_2^{k+1} \cdot GELU(W_1^{k+1} \mathbf{h}_l^{k+1} + b_1^{k+1}) + b_2^{k+1}, \\ \mathbf{z}_l^{k+1} &= LN(\mathbf{h}_l^{k+1} + \bar{\mathbf{z}}_l^{k+1}),\end{aligned}\quad (2)$$

where  $\mathbf{h}_l^{k+1}$  denotes hidden embeddings before the feed-ward operator in the  $\{k+1\}$ -th layer. Following the transformer encoder, we employ multi-head attention to jointly attend to information from different representation subspaces considering various positions, where  $W_m^{k+1}$  denotes the weight of attention heads. The input features for each layer can be used to compute three matrices:  $Q$ ,  $K$ , and  $V$  corresponding to queries, keys, and values that drive the multi-head attention block, respectively. The dot-product similarity between queries and keys determines the attention distributions of values. In detail,  $A_{l,j}^m$  denotes the attention weights between elements  $l$  and  $j$  in the  $m$ -th head, which is proportional to  $(Q_m^{k+1} \mathbf{z}_l^k)^T (K_m^{k+1} \mathbf{z}_j^k)$ .  $Q_m^{k+1}$ ,  $K_m^{k+1}$ , and  $V_m^{k+1}$  are learnable weights for the  $m$ -th attention head. Then, weight-averaged values form the output of the attention block.  $W_1^{k+1}$ ,  $W_2^{k+1}$  and  $b_1^{k+1}$ ,  $b_2^{k+1}$  are learnable weights and biases in the feed-forward layer, respectively. We use the GELU activation instead of the ReLU operator following [46]. There is a residual structure after the two sub-blocks, followed by an LN layer.

Consequently, we can acquire the embedding of the input text from the  $[CLS]$  token and the word's embedding from other tokens. Moreover, the statistical features manually extracted from the text modality, i.e.,  $\mathbf{x}^s$  containing the hashtag topic, mention, and some semantic features, such as emotional polarity [24], are widely used in the literature [24, 28] as auxiliary information. Therefore, we concatenate the embedding of the  $[CLS]$  token and the statistical features  $\mathbf{x}^s$  as the new global text embedding, i.e.,  $\hat{\mathbf{e}}_0^w = [\hat{\mathbf{e}}_0^w, \mathbf{x}^s]$ .

#### 3.3.2 Visual Encoder

In the left-bottom corner of Figure 3, we present the visual encoder  $g_v$  for extracting the image semantic embedding. The state-of-the-art Faster R-CNN [47] can obtain regional embedding and perform very well on the object detection task. Therefore, we adopt the Faster R-CNN as the visual encoder. Then, we utilize ResNet-101 to train the region proposal network (RPN), which is proven more powerful in feature extraction and performs better than VGG-16. In detail, we utilize the pre-trained Faster R-CNN to extract visual regions with pooled region of interest (ROI) embeddings for each image segment, denoted as  $\{\hat{\mathbf{e}}_l^v\}_{l=1}^{L_v}$ , where  $l$  is the index, and  $L_v$  is fixed for all image instances as [47] for a good performance. Therefore, the visual encoder can be formulated as

$$\hat{\mathbf{e}}^v = \text{Faster R-CNN}(\mathbf{x}^v), \quad (3)$$

where  $\hat{\mathbf{e}}^v$  denotes the output embeddings, and  $\mathbf{x}^v$  represents the input for Faster R-CNN. For image modality, we first resize the image to a fixed size for each instance before inputting it to the Faster R-CNN, i.e.,  $3 \times 224 \times 224$ . For training, we first use ResNet-101 to generate a feature map, which is a stack of residual blocks based on deep convolutional layers. Then, we slide a small network over the feature map to obtain a series of low-dimensional features in the RPN. These features serve as a computing set of rectangular object proposals, each with an objectiveness score. We employ a well-trained RPN to obtain proposal ROIs for regional embedding. Similar to a textual encoder, we also encode the whole image for global embedding. Considering the excellent performance of ResNet-101 in training the RPN, we employ ResNet-101 to extract the global feature embedding.

### 3.4 Objective

Thus far, we have learned the text/image global and regional embedding:  $\hat{\mathbf{e}}_0^w / \hat{\mathbf{e}}_0^v$  and  $\{\hat{\mathbf{e}}_l^w\}_{l=1}^{L_w} / \{\hat{\mathbf{e}}_l^v\}_{l=1}^{L_v}$ .  $\hat{\mathbf{e}}_0^w$  and  $\hat{\mathbf{e}}_0^v$  have different dimensions considering that  $\hat{\mathbf{e}}_0^w$  combines the statistical information, whereas  $\hat{\mathbf{e}}_l^w$  and  $\hat{\mathbf{e}}_l^v$  have the same dimensions. We aim to use these output embeddings to make the final prediction from two aspects: 1) modal aggregation and 2) intermodal inconsistency.

Modal aggregation aims to effectively combine the two modal global information for prediction. In the experiment, we made three attempts: a) embedding fusion, which adds or concatenates two modal embeddings and then inputs the fused embeddings to the FC network for prediction; b) prediction fusion, which inputs each modal global embeddings to the independent FC network for prediction and then employs the max or mean operator to fuse two modal predictions; and c) weighted fusion, which utilizes the attention mechanism on modal embeddings to calculate the weight for each modality and then combines the weight for the embedding fusion as [24] or prediction fusion as [48]. Based on the experimental verification, the embedding fusion gets the best results under the same settings on three datasets. A possible explanation is that the image contributes little to the detection, whereas the text information has more significance for detection. Therefore, the prediction or weight of the image modality may affect the final prediction more during the fusion process. Consequently, we employ embedding fusion in this paper. The modal aggregation can be formulated as

$$p_a = f(\psi_w(\hat{\mathbf{e}}_0^w) + \psi_v(\hat{\mathbf{e}}_0^v)), \quad (4)$$

where  $\psi_w/\psi_v$  is the FC layer, which projects different modal global embeddings into the same embedding space, and  $f(\cdot)$  denotes the classifier. Without any loss of generality, we utilize the FC networks here.

Intermodal inconsistency aims to measure the dissimilarity between two modalities. The direct ways include the following: 1) Global similarity (GS) – we directly use the two modal global embeddings to calculate the similarity, i.e.,  $p_c = \cos(\psi_w(\hat{\mathbf{e}}_0^w), \psi_v(\hat{\mathbf{e}}_0^v))$ . 2) Naive region similarity (NRS) – we use a naive aggregation (e.g., average or sum) for all regional embeddings to compare GS, i.e.,  $p_c = \cos(\text{ave}(\mathcal{X}^w), \text{ave}(\mathcal{X}^v))$ , where  $\text{ave}$  denotes the average operation. However, these kinds of methods mask important substructure differences. To overcome this problem, we turn to calculate the Wasserstein distance between the regional embedding of two modalities, which aims to find subtler differences rather than simple aggregations. We first illustrate the Wasserstein distance.

The Wasserstein distance [49, 50] is a function between probability distributions defined on a given metric space, which is more effective when the probability space has geometrical structures as compared with Kullback–Leibler divergences, Hellinger distance, and total variation. Intuitively, the Wasserstein distance is the minimum cost of transporting the pile of one distribution into the pile of another distribution, which formulates the problem of learning the ground metric as minimizing the difference between two polyhedral convex functions over a convex set of distance matrices. Therefore, the Wasserstein distance is powerful in such situations by considering the pairwise cost.

**Definition 1.** (Transport Polytope) For two probability vectors  $r$  and  $c$  in the simplex,  $\Gamma(r, c)$  is the transport polytope of  $r$  and  $c$ , namely, the polyhedral set of  $p \times p$  matrices:

$$\Gamma(r, c) = \{T \in \mathcal{R}_+^{p \times p} | T\mathbf{1}_p = r, T^\top \mathbf{1}_p = c\}$$

**Definition 2.** (Wasserstein Distance) Given a  $p \times p$  cost matrix  $M$ , the total cost of mapping from  $r$  to  $c$  using a transport matrix (or coupling probability)  $T$  can be quantified as  $\langle T, M \rangle$ . The Wasserstein distance problem is defined as

$$W(r, c) = \min_{T \in \Gamma(r, c)} \langle T, M \rangle$$

When  $M$  belongs to the cone of metric matrices  $\mathbb{M}$ , the value of  $W(r, c)$  is the distance [51] between  $r$  and  $c$ . In this case, assuming implicitly that  $M$  is fixed. That is, we can utilize the squared Euclidean distance of instances

for calculating  $M$ .  $M$  is a square matrix; otherwise, the corresponding position is filled with 0. Moreover, only  $r$  and  $c$  vary. We will refer to the optimal transport distance between  $r$  and  $c$ . Notably,  $W(r, c)$  is the cost of the optimal plan for transporting the predicted mass distribution  $r$  to match the target distribution  $c$ . The penalty increases when more mass is transported over longer distances based on the ground metric  $M$ . Therefore, with the region embeddings  $\mathcal{X}^v$  and  $\mathcal{X}^w$ ,  $T$  represents the node correspondences between  $\mathcal{X}^v$  and  $\mathcal{X}^w$ . We define the matrix Wasserstein distance as follows:

**Definition 3. Matrix Wasserstein Distance** Given two matrices  $\mathcal{X} \in \mathcal{R}^{m \times d}$ ,  $\mathcal{X}' \in \mathcal{R}^{n \times d}$ , in which each row represents a region embedding. The matrix Wasserstein distance can be defined as  $D_W(\mathcal{X}, \mathcal{X}') = W(\mathcal{X}, \mathcal{X}')$ .

Here, different from traditional Wasserstein distance considering continuous probability distributions, we deal with finite sets of regional embeddings. Therefore, we can reformulate the Wasserstein distance as a sum rather than an integral and use the matrix notation commonly encountered in the optimal transport [52] to represent the transportation plan. Consequently, given two sets of vectors  $\mathcal{X}^v = [\hat{\mathbf{e}}_1^v, \hat{\mathbf{e}}_2^v, \dots, \hat{\mathbf{e}}_{L_v}^v]^\top$  and  $\mathcal{X}^w = [\hat{\mathbf{e}}_1^w, \hat{\mathbf{e}}_2^w, \dots, \hat{\mathbf{e}}_{L_w}^w]^\top$ , we can define the Wasserstein distance as

$$\begin{aligned} W(\mathcal{X}^v, \mathcal{X}^w) &= \min_{T \in \Gamma(\mathcal{X}^v, \mathcal{X}^w)} \langle T, M \rangle, \\ \text{s.t. } T \mathbf{1}_{L_v} &= \mathbf{1}_{L_w} / L_w, \\ T^\top \mathbf{1}_{L_w} &= \mathbf{1}_{L_v} / L_v, \end{aligned} \quad (5)$$

where  $M$  is the distance matrix calculated by the two modalities of each instance, which contains the distances  $s(\hat{\mathbf{e}}^v, \hat{\mathbf{e}}^w)$  between each element  $\hat{\mathbf{e}}_i^v$  of  $\mathcal{X}^v$  and  $\hat{\mathbf{e}}_i^w$  of  $\mathcal{X}^w$ . We utilize the squared Euclidean distance for  $s$  here.  $T \in \Gamma$  is a transport matrix (or joint probability), and  $\langle \cdot, \cdot \rangle$  is the Frobenius dot-product. The total mass to be transported is equal to 1. Therefore, the row and column values of  $T$  must sum up to  $\mathbf{1}_{L_w} / L_w$  and  $\mathbf{1}_{L_v} / L_v$ , respectively. The transport matrix  $T$  contains the fractions that indicate how to transport the values from  $\mathcal{X}^v$  to  $\mathcal{X}^w$  with the minimal total transport effort. In summary, we can define the inconsistency prediction as

$$p_c = 1 - \mathbf{e}^{-\gamma W(\mathcal{X}^v, \mathcal{X}^w)}, \quad (6)$$

where we employ the Laplacian kernel for calculating prediction and set  $\gamma$  as 0.01 according to [53].

In summary, we can combine Eq. 4 and Eq. 6 for the final prediction:

$$L = - \sum_{i=1}^N y_i \log(\max(p_{i,c}, p_{i,a})), \quad (7)$$

where the loss function can be represented as any convex loss function, and we adopt the cross-entropy here.

The parameters in Eq. 7 include  $T$  (transport matrix),  $f$  (classifier of the modal fusion), and  $g_v/g_w$  (image/text encoder). Actually, the solution of the transport matrix  $T$  has a closed form in the forward process using  $M$ . In detail, with the calculated  $M$ ,  $T$  is the solution of an entropy-smoothed optimal transport problem:

$$T = \arg \min_T \langle T, M \rangle + \Omega(T), \quad (8)$$

where  $\Omega(T) = -\frac{1}{\lambda} \sum_{mn} T_{mn} \log T_{mn}$  is the entropy of  $T$  and  $\lambda > 0$  is the entropic regularization coefficient.  $T_{mn}$  denotes the element in the  $m$ -th row and  $n$ -th column of  $T$ . Based on the Sinkhorn theorem, we conclude that the transportation matrix can be written in the form of  $T^* = \text{diag}(u) K \text{diag}(v)$ , where  $K = \exp(-\lambda M)$  is the element-wise exponential of  $\lambda M - 1$ . Moreover, in Sinkhorn iterations,  $u$  and  $v$  are updated continuously. Taking the  $k$ -th iteration as an example, the update is represented in the following form:  $u = \frac{\mathbf{1}_{L_w} / L_w}{K^\top u^{k-1}}$  and  $v = \frac{\mathbf{1}_{L_v} / L_v}{K v^{k-1}}$ . We adopt the gradient descent to update parameters considering the loss  $L$ . The details are shown in Algorithm 1.

## 4 Experiments

In this section, extensive experiments are conducted based on three real-world datasets to demonstrate the effectiveness of our proposed method as compared with state-of-the-art methods.



**Algorithm 1** Code of the DVLFN**Input:**Data:  $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$ Parameters:  $\lambda$ **Output:**multimodal fusion network:  $f, g_v/g_w$ 


---

```

1: while stop condition is not triggered do
2:   for mini-batch sampled from  $\mathcal{D}$  do
3:     Calculate  $p_a$  according to Eq. 4;
4:     Calculate the transport matrix  $T$  according to Eq. 8;
5:     Calculate  $p_c$  according to Eq. 6;
6:     Calculate  $L$  according to Eq. 7;
7:     Update model parameters using Adam;
8:   end for
9: end while

```

---

## 4.1 Setups

### 4.1.1 Datasets

We employ three commonly used real-world rumor detection datasets: Fakeddit [25], NeuralNews [26], and Weibo [24]. The Fakeddit dataset is a multimodal dataset crawled from Reddit of fake posts. It is collected from 22 different subreddits. The samples span over almost a decade and are posted on highly active and popular pages by over 300,000 unique individual users, allowing to capture a wide variety of perspectives. The NeuralNews dataset is a constructed dataset, which consists of human and machine-generated articles with images and captions. The human-generated articles are sourced from the GoodNews dataset [54], which consists of New York Times news articles spanning from 2010 to 2018. Each news article contains a title, the main article body, and image-caption pairs. For fairness, we remove the image caption to keep the data format consistent with that in Fakeddit. The Weibo dataset is specifically constructed for rumor detection. In detail, considering the size and timeliness of the Weibo dataset in [24], we extend the original dataset by crawling more data from January 2018 to February 2019 with the same technique as [24]. Therefore, the Weibo dataset is an expanded dataset of the original paper. In this dataset, non-rumors are collected from authoritative news sources in China, such as Xinhua News Agency. The rumors are crawled from the website and verified by the official rumor debunking system of Weibo. This system encourages common users to report suspicious posts and examines suspicious posts by a committee of trusted users. We follow the same steps in the work [24]. We first remove the duplicated and low-quality images to ensure the quality of the entire dataset. Then, we apply a single-pass clustering method [55] to discover newly emerged events from posts. Finally, we split the whole dataset into training, validation, and testing sets in a 7:1:2 ratio and ensure that they do not contain any common event. The descriptions of the training and testing data distributions are listed in Table 1.

**Table 1** Statistics of three datasets.

Statistics		Fakeddit	NeuralNews	Weibo
Training set	Rumor	222081	22400	11792
	Non-rumor	341919	22400	11874
Validation set	Rumor	23320	3200	1678
	Non-rumor	36022	3200	1702
Testing set	Rumor	23507	6400	3371
	Non-rumor	35812	6400	3389
All		682661	64000	33806

**Table 2** Classification results of different methods on three datasets. The best results are highlighted in bold. “Acc,” “Pre,” and “Rec” denote the accuracy, precision, and recall.

Methods	Fakeddit				NeuralNews				Weibo			
	Acc	Pre	Rec	F1	Acc	Pre	Rec	F1	Acc	Pre	Rec	F1
Textual	.868	.895	.886	.890	.930	.920	.942	.931	.959	.956	.963	.959
Visual	.794	.842	.811	.826	-	-	-	-	.847	.851	.837	.844
Statistical	.810	.838	.851	.844	.819	.820	.816	.818	.909	.907	.924	.891
Mean-LF	.891	.907	.914	.910	.930	.933	.927	.930	.955	.954	.955	.955
Max-LF	.892	.908	.912	.910	.931	.932	.930	.931	.945	.952	.936	.944
Weighted-LF	.890	.903	.917	.910	.924	.922	.928	.925	.947	.961	.932	.946
Add-EF	.901	.920	.916	.918	.933	.940	.926	.933	.962	.969	.954	.961
Concatenate-EF	.897	.919	.911	.915	.930	.934	.925	.929	.961	.964	.957	.960
Weighted-EF	.895	.908	.919	.914	.930	.943	.915	.928	.960	.971	.948	.959
GAN-GRU	.866	.881	.899	.890	.892	.877	.913	.895	.964	.968	.959	.963
MVNN	.768	.828	.778	.802	-	-	-	-	.842	.849	.826	.838
att-RNN	.901	.916	.921	.918	.893	.881	.909	.895	.964	.972	.955	.964
EANN	.871	.898	.887	.892	.857	.872	.836	.854	.953	.966	.939	.952
MAVE	.884	.911	.896	.903	.891	.884	.901	.892	.964	.980	.948	.964
VL-BERT	.888	.909	.907	.908	.921	.941	.898	.919	.946	.957	.934	.945
DRMM	.906	.929	.915	.922	.947	.967	.927	.946	.968	.968	.966	.967
MKN	.884	.898	.912	.905	.892	.906	.872	.889	.959	.970	.948	.959
SpotFake	.897	.915	.915	.915	.932	.930	.935	.932	.959	.967	.953	.960
CARMN	.888	.913	.900	.906	.891	.893	.887	.890	.966	.967	.965	.966
DIDAN	.893	.919	.903	.911	.956	.952	.960	.956	.967	.964	.969	.967
SAFE	.886	.912	.897	.905	.884	.890	.876	.883	.963	.974	.953	.963
EM-FEND	.905	.925	.917	.921	.951	.967	.934	.950	.967	.970	.966	.968
DVLFN	<b>.915</b>	<b>.931</b>	<b>.928</b>	<b>.929</b>	<b>.982</b>	<b>.989</b>	<b>.975</b>	<b>.982</b>	<b>.978</b>	<b>.981</b>	<b>.974</b>	<b>.977</b>

#### 4.1.2 Implementation

For the statistical features, we take the most commonly used social features in the literature following [24], which extracts 12-dimensional statistical features. The image encoder employs the Faster R-CNN [47], which uses fixed 36 ROIs. Each ROI and the global embedding are learned using ResNet-101 [56]. The text encoder employs the transformer [8], which has 12 layers of transformer blocks, where each block has 12 self-attention heads. The parameters are initialized from the BERT base, which is pre-trained on text data only. The model is trained with the following hyperparameters: ADAM optimizer [57] with a learning rate of 0.001, epoch number of 50, mini-batch size of 64, and weight decay of 0.001. The output layer size of each model is 200; i.e., each modal’s final embedding is 200-dimensional. We implement the model using the MindSpore Lite tool <sup>3)</sup>.

#### 4.1.3 Comparison of Methods

To validate the proposed model on the rumor detection task, we compare our model with four groups of baseline methods, which are widely employed in the literature: 1) Single-modal methods: textual, visual, statistical, GAN-GRU [19], and MVNN [31]. 2) Direct multimodal early or late fusion methods: early fusion, late fusion, and weighted fusion. 3) Deep multimodal fusion methods: att-RNN [24], EANN [23], MVAE [22], VL-BERT [58], DRMM [59], MKN [60], SpotFake [61], and CARMN [62]. 4) Deep multimodal fusion method with inconsistency measure: DIDAN [26], SAFE [41] and EM-FEND [42].

<sup>3)</sup> <https://www.mindspore.cn/2020>

**Table 3** T-test results of different methods on three datasets.

Methods	Fakeddit				NeuralNews				Weibo			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
Add-EF	0.001	0.013	0.004	0.000	0.000	0.000	0.000	0.000	0.000	0.010	0.004	0.000
att-RNN	0.000	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.006	0.001	0.000
EANN	0.000	0.000	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.002	0.000	0.000
MVAE	0.000	0.006	0.002	0.000	0.000	0.000	0.000	0.000	0.000	0.004	0.008	0.000
DIDAN	0.000	0.003	0.001	0.000	0.000	0.001	0.000	0.000	0.002	0.013	0.001	0.001
DRMM	0.000	0.017	0.012	0.000	0.000	0.003	0.000	0.000	0.000	0.009	0.002	0.001
MKN	0.000	0.001	0.006	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.003	0.000
SAFE	0.000	0.001	0.002	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.000
SpotFake	0.000	0.002	0.000	0.000	0.000	0.000	0.001	0.000	0.001	0.002	0.000	0.000
CARMN	0.000	0.001	0.002	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.004	0.000
EM-FEND	0.001	0.012	0.007	0.001	0.000	0.002	0.001	0.000	0.002	0.015	0.003	0.000

- **Textual:** A single deep model with texts only. We use BERT in this study.
- **Visual:** A single deep model with images only. We use ResNet101 in this study.
- **Statistical:** A single tree model with statistical features only. We use the SOTA decision tree model, i.e., LightGBM [63].
- **GAN-GRU:** A GAN-style approach. A generator is designed to produce augmented posts as rumors to pressurize the discriminator to learn strong rumor-indicative representations.
- **MVNN:** A deep approach that fuses the visual information of frequency and pixel domains for detecting fake news.
- **Early Fusion:** Direct deep multimodal early fusion method. We employ visual and textual transformers as basic models and then add or concatenate different modal embeddings for final prediction, which is denoted as Add-EF/Concatenate-EF for simplicity.
- **Late Fusion:** Direct deep multimodal late fusion method. We train each modal deep model independently and then utilize the mean/max operator for the final predictions, which is denoted as Mean-LF/MAX-LF for simplicity.
- **Weighted Fusion:** Direct deep multimodal fusion method. We train the weight network to learn weights for the image and text and then utilize the weighted sum for modal embeddings or final predictions, which is denoted as Weighted-EF/Weighted-LF for simplicity.
- **VL-BERT:** A single flow-based multimodal transformer model, which concatenates the image-text input for joint training.
- **att-RNN:** A deep multimodal fusion network, which consists of an innovative RNN and an attention mechanism for fusing textual, visual, and statistical features.
- **EANN:** A deep multimodal fusion network, which employs a multitask loss, including rumor classification and event discrimination for learning common embeddings.
- **MVAE:** A deep multimodal fusion network, which trains a multimodal variational autoencoder jointly with a rumor detector.
- **DRMM:** A deep multimodal fusion network, which conducts deep interactions between images and sentences for modality feature aggregation.
- **DIDAN:** A deep multimodal fusion network considering inconsistencies, which manually designs a binary named entity indicator as an extra input.
- **MKN:** A deep multimodal knowledge-aware network considering multimodal content and external knowledge-level connections.
- **SAFE:** A deep multimodal fusion network considering cross-modal inconsistency, which computes multimodal relevance as the auxiliary objective for fake news detection.
- **SpotFake:** A deep multimodal fusion network that concatenates the textual and visual features obtained from pre-trained models.
- **CARMN:** A deep multimodal fusion network, which employs a cross-modal attention residual network for fusing multimodal features.

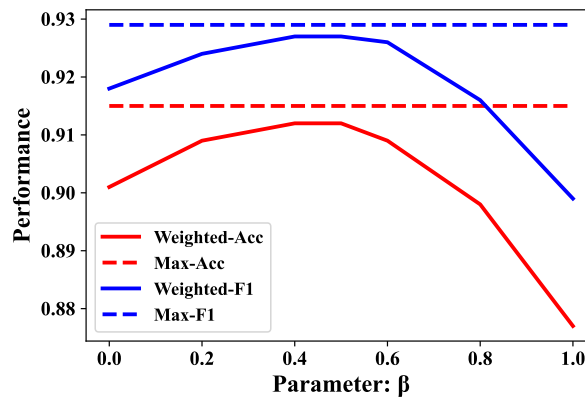
**Table 4** Ablation study of different encoders on three datasets. The best results are highlighted in bold.

Text	Image	Fakeddit				NeuralNews				Weibo			
		Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
BERT	-	.868	.895	.886	.890	.930	.920	.942	.931	.959	.956	.963	.959
LSTM	-	.865	.887	.890	.889	.883	.907	.852	.879	.956	.950	.961	.956
GRU	-	.865	.885	.892	.888	.891	.876	.911	.893	.957	.961	.953	.957
-	ResNet101	.794	.842	.811	.826	-	-	-	-	.847	.851	.837	.844
-	Vgg16	.766	.809	.801	.805	-	-	-	-	.840	.882	.779	.827
-	EfficientNet	.738	.788	.775	.782	-	-	-	-	.834	.831	.831	.831
BERT	ResNet101	.906	.921	.924	.922	.938	.942	.934	.938	.963	.960	.966	.963
BERT	Vgg16	.902	.923	.915	.919	.937	.937	.937	.937	.960	.964	.956	.960
BERT	EfficientNet	.899	.914	.920	.917	.932	.938	.925	.931	.955	.959	.949	.954
LSTM	ResNet101	.903	.919	.921	.920	.890	.891	.888	.889	.957	.961	.953	.957
LSTM	Vgg16	.895	.916	.909	.913	.894	.901	.885	.893	.954	.956	.952	.954
LSTM	EfficientNet	.896	.910	.918	.914	.890	.899	.879	.889	.953	.965	.939	.952
GRU	ResNet101	.897	.916	.912	.914	.894	.895	.894	.894	.961	.973	.947	.960
GRU	Vgg16	.890	.916	.900	.908	.896	.901	.890	.895	.960	.963	.956	.959
GRU	EfficientNet	.884	.901	.908	.904	.894	.909	.877	.892	.954	.972	.934	.953
DVLFN		<b>.915</b>	<b>.931</b>	<b>.928</b>	<b>.929</b>	<b>.982</b>	<b>.989</b>	<b>.975</b>	<b>.982</b>	<b>.978</b>	<b>.981</b>	<b>.974</b>	<b>.977</b>

• **EM-FEND**: A deep multimodal fusion network considering cross-modal inconsistency, which models the complementary text information, mutual enhancement, and entity inconsistency.

To further investigate the impact of each item in the proposed model, we design several ablation studies for comparison: 1) **w/o prediction**: The variant of the DVLFN that removes the prediction term. 2) **w/o inconsistency**: The variant of the DVLFN that removes the inconsistency term. 3) **w/o statistical**: The variant of the DVLFN that removes the statistical information. 4) **Different Encoders**: We also design various variants of the DVLFN with different visual/textual basic encoders, i.e., BERT+ResNet101, BERT+Vgg16, BERT+EfficientNet, LSTM+ResNet101, LSTM+Vgg16, LSTM+EfficientNet, GRU+ResNet101, GRU+Vgg16, and GRU+EfficientNet.

Considering that existing approaches have already mentioned the disadvantages of linear methods [22–24], we do not include linear methods in this study.

**Figure 4** Performance of the DVLFN with different  $\beta$  for weighted prediction.

**Table 5** Ablation studies of the DVLFN and variants on three datasets. The best results are highlighted in bold.

Methods	Fakeddit				NeuralNews				Weibo			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
w/o statistical	.910	.923	.926	.926	.980	.987	.970	.980	.972	.968	.973	.972
w/o prediction	.877	.896	.901	.899	.971	.984	.957	.970	.929	.949	.906	.927
w/o inconsistency	.901	.920	.916	.918	.933	.940	.926	.933	.962	.969	.954	.961
GS	.906	.921	.924	.922	.938	.942	.934	.938	.963	.960	.966	.963
NRS	.906	.928	.916	.922	.946	.928	.966	.947	.965	.977	.952	.964
Hungarian	.911	.926	.927	.927	.976	.983	.970	.976	.964	.970	.957	.963
DVLFN	<b>.915</b>	<b>.931</b>	<b>.928</b>	<b>.929</b>	<b>.982</b>	<b>.989</b>	<b>.975</b>	<b>.982</b>	<b>.978</b>	<b>.981</b>	<b>.974</b>	<b>.977</b>



Argetina's football squad celebrate a goal during the 2018 World Cup.



A photo of a horse wearing horse shoes before the metal variety were invented.



Mehmet II celebrates the capture of Constantinople with his leading generals.



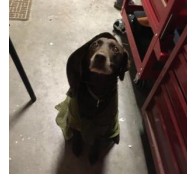
Report: Majority Of Americans Now Eating One Continuous Meal A Day.



First single-celled organisms perform mitosis.



British medic risks his life to save an injured comrade at the Battle of the Somme.

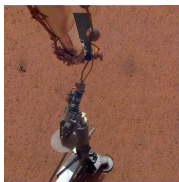


A young Marine waiting for orders during Desert Storm.



Al Qaeda terrorists first attempt at hitting the twin towers using kites and helium balloons.

(a) Examples of Rumors



Robot 'Mole' on Mars Begins Digging Into Red Planet This Week.



Piece of bark that looks like a dog/giraffe.



Community college launches free food pantry initiative on campus to address food insecurity issues for students.



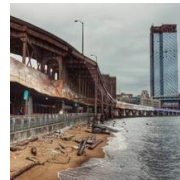
This small rainbow during a storm.



A water spot on the road in the shape of a perfect heart.



Olympian sells medal to help pay for child's cancer treatment.



A stretch of beach in southern Manhattan - under the Brooklyn Bridge looking east.



This exhibit illustrates the effects of water erosion over time with rock.

(b) Examples of Non-rumors

**Figure 5** (Best viewed when zoomed in.) Several top rumors and non-rumors detected by the DVLFN from the Fakeddit dataset.

## 4.2 Classification Results

In this section, we present the comparison results of the baselines and ablation methods with four criteria, i.e., accuracy, precision, recall, and F1 score. We only show the best results as [22, 24].

### 4.2.1 Comparison with the Baseline Methods

Table 2 records the test results of all compared models. The results in the visual model are represented as “-” because there are no fake images in NeuralNews. NeuralNews is a constructed dataset that uses GLOVER [64] to generate fake articles with original titles and articles, and the raw images remain unchanged.

We draw the following observations: 1) The performance of texts is better than that of images on all datasets considering various criteria because texts always contain more semantic information for classification. Moreover, the distinct performance indirectly explains the inconsistency between images and texts. 2) The classification results of statistical information are also good, which validates the effectiveness of the statistical modality. 3) Among the different fusion methods, the performance of early fusion with the add operator performs the best on all the datasets considering various criteria. Therefore, we adopt it in the modal aggregation module. 4) The GAN-based model, i.e., GAN-GRU, improves slightly as compared with the original model, i.e., GRU. 5) The EANN and MVAE are not even as good as the single-modal models on partial settings because of the following reasons: a) The EANN needs extra supervision information for multitask training, which is missing on both datasets. b) The MVAE maps images and texts to the same vector space for a consistent constraint, but it brings confusion for embedding learning on the dataset (i.e., NeuralNews) with a large number of inconsistent instances. Hence, the MVAE performs worse on the NeuralNews dataset. 6) VL-BERT performs worse than independent training approaches (e.g., DVLFN, late fusion, and early fusion), which validates the advantage of the independent encoder for each modality. 7) The SOTA deep multimodal fusion method, i.e., DRMM, performs better than the single-modal and direct fusion models because they adopt a sophisticated integration strategy. 8) The fusion model considering inconsistency, i.e., DIDAN, performs well on the constructed dataset (i.e., NeuralNews), but it performs worse on real datasets, which indicates that the generalization of artificially constructed inconsistent features is not good. 9) The DVLFN performs better than the DIDAN and EM-FEND, which validates that the adaptively learned inconsistency is even better. The EM-FEND performs better than the DIDAN on most settings because the EM-FEND considers more OCR information. 10) The DVLFN outperforms other baseline methods in all cases. We attribute the superiority to the intermodal inconsistency measurement. In addition, considering the difficulty of different datasets, the DVLFN promotes limited effects on the Fakeddit dataset, i.e., approximately 0.7% on F1 score, and it significantly improves on NeuralNews, i.e., approximately 2.7%. 11) To validate the significance and prove the superiority of the DVLFN, we perform statistical significance tests. Table 3 records the results of comparing state-of-the-art methods. We find that our proposed DVLFN indeed outperforms other algorithms significantly. For example, the p-values of the DVLFN are always lower than 0.05. 12) We conduct more experiments on the raw Weibo dataset [24]. Considering the F1 score, the best comparison method, i.e., EM-FEND, acquires 0.901, and the DVLFN gets 0.919, which also validates the effectiveness of the DVLFN on small datasets.

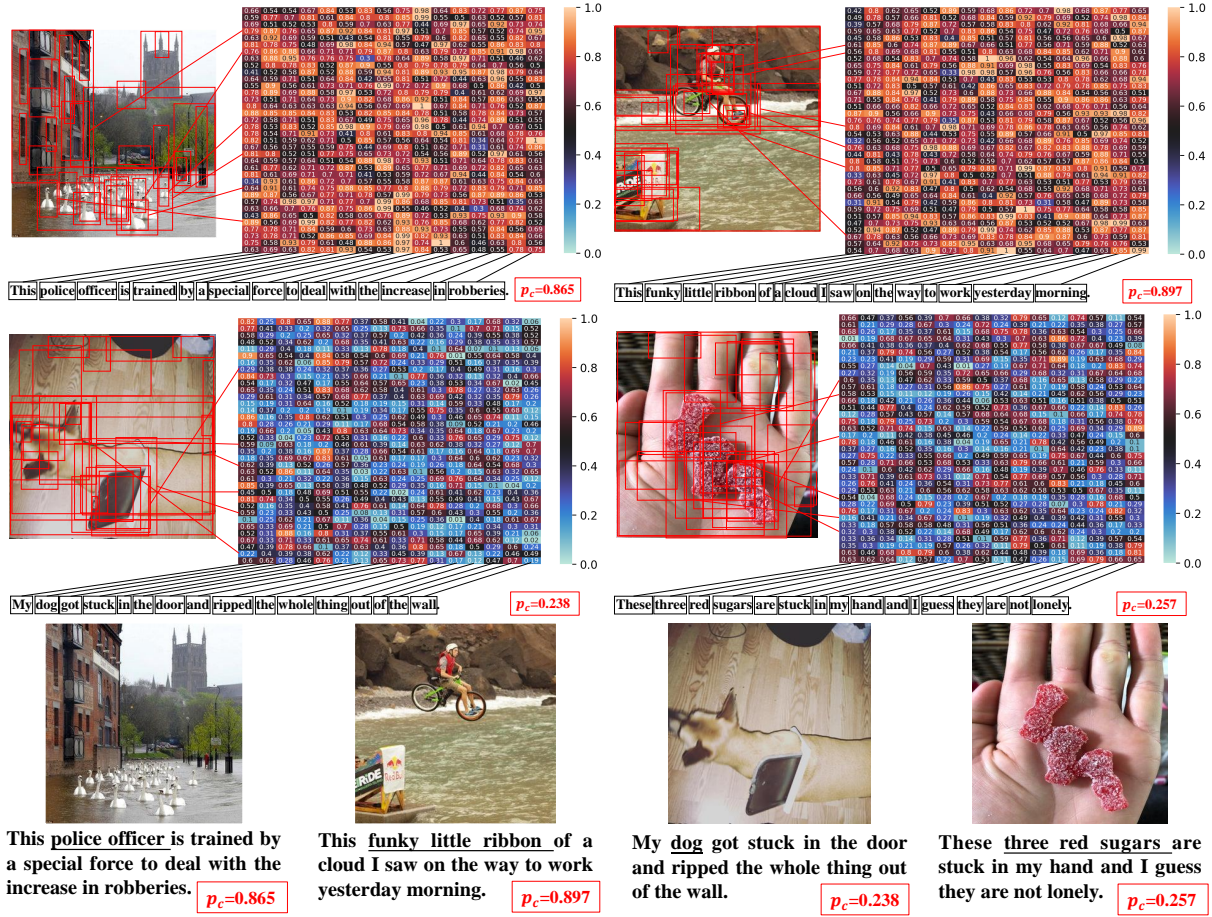
### 4.2.2 Ablation Study

**Study of variants of the DVLFN.** The results are shown in Table 5, from which we can conclude the following: 1) Using the inconsistency measure alone, i.e., w/o prediction, can improve the results. For example, the w/o prediction method achieves significant performance on the NeuralNews dataset, considering that NeuralNews is a constructed inconsistent dataset. 2) The w/o inconsistency method performs well on the three datasets, which validates the selection of the encoder. 3) The w/o statistical variant method is worse than the DVLFN but better than the w/o prediction and w/o inconsistency methods. This finding reveals that the statistical information is helpful to the prediction, but the effect is small.

**Study on different encoders.** To explore the effectiveness of different deep encoders for different modalities, we conduct more ablation studies. Related results are recorded in Table 4. For image and multimodal classifications, ResNet101 performed the best, followed by the other two encoders. Moreover, BERT achieved better results than the LSTM and GRU for text modality.

**Study on each loss term.** To explore the importance of the fusion term and inconsistency term, we adopt the weighted prediction replacing the max operation, i.e.,  $L = -\sum_{i=1}^N y_i \log((1-\beta)p_{i,a} + \beta p_{i,c}) + (1-y_i) \log(1 - ((1-\beta)p_{i,a} + \beta p_{i,c}))$ . Figure 4 records the results of the weighted prediction (i.e., tuning the parameter  $\beta$  in  $\{0, 0.2, 0.4, 0.6, 0.8, 1\}$ ) on the Fakeddit dataset. The results reveal the following: 1) Excessive  $\beta$  will affect the performance. The reason is that content-based predictions are more important, and the inconsistency measure tends





**Figure 6** Examples with high inconsistency and low inconsistency scores detected by the DVLFN on the NeuralNews dataset. The two cases in the first row are multimodal examples with high inconsistency scores, and the two cases in the second row are multimodal examples with low inconsistency scores. The bottom row provides the raw image-text pairs for viewing convenience. The normalized matrix represents the distance matrix  $M$  calculated by image regions and text words during the computation of inconsistencies.  $p_c$  denotes the inconsistency score.

to detect rumors whose contents are difficult to classify. However, the modalities are obviously inconsistent. 2) The max operation achieves the best results because only when the modalities are extremely inconsistent will it be considered a rumor. Thus, the max operation ensures a better prediction.

**Study on different distance measures.** We conducted more studies to validate the role of intermodal inconsistency: 1) Different distance measures, i.e., GS and NRS. 2) Different graph matching methods, i.e., Wasserstein (i.e., DVLFN) and Hungarian. Table 5 records the results with different distance measures. The DVLFN achieves the best performance on various criteria of three datasets, which further verifies the important role of considering region-based intermodal inconsistency.

### 4.3 Case Study

We provide several successful examples for the qualitative analysis of the DVLFN. In detail, we sort the detected rumors and non-rumors based on the prediction scores of DVLFN and illustrate several demos from each class of the Fakeddit dataset in Figure 5. The portraits reveal that the DVLFN takes advantage of textual and visual content for rumor detection. For example, rumor examples carry significant fake patterns in text and image contents as compared with non-rumor ones. Furthermore, to illustrate the importance of inconsistency, we add cases to locate inconsistencies in images and texts according to the similarity matrix. Figure 6 indicates the results. The matrix represents the distance matrix  $M$  calculated by the image regions and text words during the computation of inconsistencies. The values in the matrix are the normalized distances between image regions and text words, where large values denote dissimilarity.  $p_c$  denotes the inconsistency score, where small values represent high consistency. The rumor cases with high inconsistency predictions are always with large regional distances. For example, in the first case in the top row of Figure 6, the image describes objects of “goose,” “water,” and “building,” whereas the text describes objects of “police” and “robberies.” The contents between the images and texts are inconsistent, and

the distance matrix can effectively discover regions with large inconsistency values. On the contrary, consistent image–text pairs (i.e., the examples in the second row of Figure 6) can well locate the consistency of images and texts.

## 5 Conclusions

Rumors on social media have posed great challenges to supervisors. However, most of the current multimodal fusion methods ignored a notable characteristic of rumors, i.e., inconsistency between modalities. In this study, we develop the novel DVLFN, which detects rumors with the aid of a fine-grained inconsistency measure. In detail, the DVLFN first utilizes visual–linguistic deep embeddings to calculate regional embeddings, which aims to measure intermodal inconsistencies. Then, the DVLFN composes a reliable fusion mechanism to process concatenation for the final prediction. Experiments are conducted on real-world multimedia rumor detection datasets, and the results show the superior performance of the proposed DVLFN.

**Acknowledgements** This research was supported by the NSFC (62006118, 61906092, 61773198, 91746301), Natural Science Foundation of Jiangsu Province of China under Grant (BK20200460, BK20190441), Jiangsu Shuangchuang (Mass Innovation and Entrepreneurship) Talent Program, and CAAI-Huawei MindSpore Open Fund (CAAIXSJLJ-2021-014B).

## References

- Gordon W Allport and Leo Postman. The psychology of rumor. 1947.
- Hunt Allcott and Matthew Gentzkow. Social media and fake news in the 2016 election. *Journal of economic perspectives*, 31(2):211–36, 2017.
- Ceren Budak. What happened? the spread of fake news publisher content during the 2016 U.S. presidential election. In *Proceedings of the The World Wide Web Conference*, pages 139–150, San Francisco, CA, 2019.
- Clement Farabet, Camille Couprie, Laurent Najman, and Yann LeCun. Learning hierarchical features for scene labeling. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1915–1929, 2013.
- Yang Yang, De-Chuan Zhan, Yi-Feng Wu, Zhi-Bin Liu, Hui Xiong, and Yuan Jiang. Semi-supervised multi-modal clustering and classification with incomplete modalities. *IEEE Trans. Knowl. Data Eng.*, 33(2):682–695, 2021.
- Ronan Collobert and Jason Weston. A unified architecture for natural language processing: deep neural networks with multitask learning. In *Proceedings of the 25th International Conference Machine Learning*, pages 160–167, Helsinki, Finland, 2008.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *Proceedings of the 3rd International Conference on Learning Representations*, San Diego, CA, 2015.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems 30*, pages 5998–6008, Long Beach, CA, 2017.
- Manish Gupta, Peixiang Zhao, and Jiawei Han. Evaluating event credibility on twitter. In *Proceedings of the SIAM International Conference on Data Mining*, pages 153–164, Anaheim, California, 2012.
- Sejeong Kwon, Meeyoung Cha, Kyomin Jung, Wei Chen, and Yajun Wang. Prominent features of rumor propagation in online social media. In *Proceedings of the IEEE 13th International Conference on Data Mining*, pages 1103–1108, Dallas, TX.
- Ke Wu, Song Yang, and Kenny Q. Zhu. False rumors detection on sina weibo by propagation structures. In *Proceedings of the IEEE International Conference on Data Engineering*, pages 651–662, Seoul, South Korea, 2015.
- Zhiwei Jin, Juan Cao, Yongdong Zhang, and Jiebo Luo. News verification by exploiting conflicting social viewpoints in microblogs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2972–2978, Phoenix, Arizona, 2016.
- Jing Ma, Wei Gao, Prasenjit Mitra, Sejeong Kwon, Bernard J. Jansen, Kam-Fai Wong, and Meeyoung Cha. Detecting rumors from microblogs with recurrent neural networks. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pages 3818–3824, New York, NY, 2016.
- Feng Yu, Qiang Liu, Shu Wu, Liang Wang, and Tieniu Tan. A convolutional approach for misinformation identification. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 3901–3907, Melbourne, Australia, 2017.
- Christina Boididou, Symeon Papadopoulos, Duc-Tien Dang-Nguyen, Giulia Boato, and Yiannis Kompatsiaris. The certh-unitn participation@ verifying multimedia use 2015. In *MediaEval*, 2015.
- Peng Qi, Juan Cao, Tianyun Yang, Junbo Guo, and Jintao Li. Exploiting multi-domain visual information for fake news detection. In *Proceedings of the IEEE International Conference on Data Mining*, pages 518–527, Beijing, China, 2019.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27*, pages 2672–2680, Quebec, Canada, 2014.
- Lakshmanan Nataraj, Tajuddin Manhar Mohammed, B. S. Manjunath, Shivkumar Chandrasekaran, Arjuna Flenner, Jawadul H. Bappy, and Amit K. Roy-Chowdhury. Detecting GAN generated fake images using co-occurrence matrices. In *Proceedings of the Media Watermarking, Security, and Forensics*, Burlingame, CA, 2019.
- Jing Ma, Wei Gao, and Kam-Fai Wong. Detect rumors on twitter by promoting information campaigns with generative adversarial learning. In *Proceedings of the World Wide Web Conference*, pages 3049–3055, San Francisco, CA, 2019.
- Bin-Bin Jia and Min-Ling Zhang. Multi-dimensional classification via selective feature augmentation. *Machine Intelligence Research*, 19(1):38–51, 2022.
- Huaiwen Zhang, Quan Fang, Shengsheng Qian, and Changsheng Xu. Multi-modal knowledge-aware event memory network for social media rumor detection. In *Proceedings of the ACM International Conference on Multimedia*, pages 1942–1951, Nice, France, 2019.
- Dhruv Khattar, Jaipal Singh Goud, Manish Gupta, and Vasudeva Varma. MVAE: multimodal variational autoencoder for fake news detection. In *Proceedings of the World Wide Web Conference*, pages 2915–2921, San Francisco, CA, 2019.
- Yaqing Wang, Fenglong Ma, Zhiwei Jin, Ye Yuan, Guangxu Xun, Kishlay Jha, Lu Su, and Jing Gao. EANN: event adversarial neural networks for multi-modal fake news detection. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 849–857, London, UK, 2018.
- Zhiwei Jin, Juan Cao, Han Guo, Yongdong Zhang, and Jiebo Luo. Multimodal fusion with recurrent neural networks for rumor detection on microblogs. In *Proceedings of the ACM on Multimedia Conference*, pages 795–816, Mountain View, CA, 2017.
- K. Nakamura, S. Levy, and W. Y. Wang. r/fakeddit: A new multimodal benchmark dataset for fine-grained fake news detection. 2019.
- Reuben Tan, Bryan A. Plummer, and Kate Saenko. Detecting cross-modal inconsistency to defend against neural fake news. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 2081–2106, 2020.



- 27 Zhiwei Jin, Juan Cao, Yu-Gang Jiang, and Yongdong Zhang. News credibility evaluation on microblog with a hierarchical propagation model. In *Proceedings of the IEEE International Conference on Data Mining*, pages 230–239, Shenzhen, China, 2014.
- 28 Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. Information credibility on twitter. In *Proceedings of the International Conference on World Wide Web*, pages 675–684, Hyderabad, India, 2011.
- 29 Zhiwei Jin, Juan Cao, Yongdong Zhang, Jianshe Zhou, and Qi Tian. Novel visual and statistical image features for microblogs news verification. *IEEE transactions on multimedia*, 19(3):598–608, 2016.
- 30 Han Guo, Juan Cao, Yazhi Zhang, Junbo Guo, and Jintao Li. Rumor detection with hierarchical social attention network. In *Proceedings of the ACM International Conference on Information and Knowledge Management*, pages 943–951, Torino, Italy, 2018.
- 31 Peng Qi, Juan Cao, Tianyun Yang, Junbo Guo, and Jintao Li. Exploiting multi-domain visual information for fake news detection. In *Proceedings of the IEEE International Conference on Data Mining*, pages 518–527, Beijing, China, 2019.
- 32 Christina Boididou, Katerina Andreadou, Symeon Papadopoulos, Duc-Tien Dang-Nguyen, Giulia Boato, Michael Riegler, and Yiannis Kompatsiaris. Verifying multimedia use at mediaeval 2015. In *Working Notes Proceedings of the MediaEval 2015 Workshop*, Wurzen, Germany, 2015.
- 33 Andrej Karpathy and Fei-Fei Li. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015.
- 34 Yang Yang, Yi-Feng Wu, De-Chuan Zhan, Zhi-Bin Liu, and Yuan Jiang. Deep robust unsupervised multi-modal network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5652–5659, Honolulu, Hawaii, 2019.
- 35 Yang Yang, Chubing Zhang, Yi-Chu Xu, Dianhai Yu, De-Chuan Zhan, and Jian Yang. Rethinking label-wise cross-modal retrieval from A semantic sharing perspective. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pages 3300–3306, Virtual, 2021.
- 36 Qi Wu, Damien Teney, Peng Wang, Chunhua Shen, Anthony R. Dick, and Anton van den Hengel. Visual question answering: A survey of methods and datasets. *Comput. Vis. Image Underst.*, 163:21–40, 2017.
- 37 Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6077–6086, Salt Lake City, UT, 2018.
- 38 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. *CoRR*, abs/2102.05918, 2021.
- 39 Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollar, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the IEEE European Conference on Computer Vision*, pages 740–755, Zurich, Switzerland, 2014.
- 40 Mark J. Huiskes and Michael S. Lew. The MIR flickr retrieval evaluation. In *Proceedings of the ACM International Conference on Multimedia*, pages 39–43, British Columbia, Canada, 2008.
- 41 Xinyi Zhou, Jindi Wu, and Reza Zafarani. SAFE: similarity-aware multi-modal fake news detection. In *PAKDD* (2), pages 354–367, Singapore, 2020.
- 42 Peng Qi, Juan Cao, Xirong Li, Huan Liu, Qiang Sheng, Xiaoyue Mi, Qin He, Yongbiao Lv, Chenyang Guo, and Yingchao Yu. Improving fake news detection by using an entity-enhanced framework to fuse diverse multimodal clues. In *ACM Multimedia*, pages 1212–1220, Virtual Event, China, 2021.
- 43 Yang Yang, Han-Jia Ye, De-Chuan Zhan, and Yuan Jiang. Auxiliary information regularized machine for multiple modality feature learning. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence*, pages 1033–1039, Buenos Aires, Argentina, 2015.
- 44 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4171–4186, Minneapolis, MN, 2019.
- 45 Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144, 2016.
- 46 Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian error linear units. *CoRR*, abs/1606.08415, 2016.
- 47 Kuang-Huei Lee, Xi Chen, Gang Hua, Houdong Hu, and Xiaodong He. Stacked cross attention for image-text matching. In *Proceedings of the European Conference Computer Vision*, pages 212–228, Munich, Germany, 2018.
- 48 Yang Yang, Ke-Tao Wang, De-Chuan Zhan, Hui Xiong, and Yuan Jiang. Comprehensive semi-supervised multi-modal learning. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pages 4092–4098, Macao, China, 2019.
- 49 R Yossi, LJ Guibas, and C Tomasi. The earth mover’s distance multi-dimensional scaling and color-based image retrieval. In *ARPA*, 1997.
- 50 Yang Yang, Zhao-Yang Fu, De-Chuan Zhan, Zhi-Bin Liu, and Yuan Jiang. Semi-supervised multi-modal multi-instance multi-label deep network with optimal transport. *IEEE Trans. Knowl. Data Eng.*, 33(2):696–709, 2021.
- 51 Cedric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- 52 Yossi Rubner, Carlo Tomasi, and Leonidas J. Guibas. The earth mover’s distance as a metric for image retrieval. *Int. J. Comput. Vis.*, 40(2):99–121, 2000.
- 53 Matteo Togninalli, M. Elisabetta Ghisu, Felipe Llinares-López, Bastian Rieck, and Karsten M. Borgwardt. Wasserstein weisfeiler-lehman graph kernels. In *Advances in Neural Information Processing Systems* 32, pages 6436–6446, Vancouver, Canada, 2019.
- 54 Ali Furkan Biten, Lluís Gómez, Marçal Rusiñol, and Dimosthenis Karatzas. Good news, everyone! context driven entity-aware captioning for news images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12466–12475, Long Beach, CA, 2019.
- 55 Zhiwei Jin, Juan Cao, Yu-Gang Jiang, and Yongdong Zhang. News credibility evaluation on microblog with a hierarchical propagation model. In *Proceedings of the IEEE International Conference on Data Mining*, pages 230–239, Shenzhen, China, 2014.
- 56 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, Las Vegas, NV, 2016.
- 57 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations*, San Diego, CA, 2015.
- 58 Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. VL-BERT: pre-training of generic visual-linguistic representations. In *Proceedings of the International Conference on Learning Representations*, Addis Ababa, Ethiopia, 2020.
- 59 Meihan Tong, Shuai Wang, Yixin Cao, Bin Xu, Juanzi Li, Lei Hou, and Tat-Seng Chua. Image enhanced event detection in news articles. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, pages 9040–9047, New York, NY, 2020.
- 60 Huaiwen Zhang, Quan Fang, Shengsheng Qian, and Changsheng Xu. Multi-modal knowledge-aware event memory network for social media rumor detection. In *ACM Multimedia*, pages 1942–1951, Nice, France, 2019.
- 61 Shivangi Singhal, Rajiv Ratn Shah, Tanmoy Chakraborty, Ponnurangam Kumaraguru, and Shin’ichi Satoh. Spotfake: A multi-modal framework for fake news detection. In *BigMM*, pages 39–47, Singapore, 2019.
- 62 Chenguang Song, Nianwen Ning, Yunlei Zhang, and Bin Wu. A multimodal fake news detection model based on crossmodal attention residual and multichannel convolutional neural networks. *Inf. Process. Manag.*, 58(1):102437, 2021.
- 63 Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems* 30, pages 3146–3154, Long Beach, CA, 2017.
- 64 Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. Defending against neural fake news. In

*Advances in Neural Information Processing Systems* 32, pages 9051–9062, Vancouver, Canada, 2019.