MViR: Multi-View Visual-Semantic Representation for Fake News Detection

Anonymous ACL submission

Abstract

With the rise of online social networks, detecting fake news accurately is essential for a healthy online environment. While existing methods have advanced multimodal fake news detection, they often neglect the multi-view visual-semantic aspects of news, such as different text perspectives of the same image. To address this, we propose a Multi-View Visual-Semantic Representation (MViR) framework. Our approach includes a Multi-View Representation module using pyramid dilated convolution to capture multi-view visual-semantic features, a Multi-View Feature Fusion module to integrate these features with text, and multiple aggregators to extract multi-view semantic cues for detection. Experiments on benchmark datasets demonstrate the superiority of MViR. The codes will be released.

1 Introduction

011

012

013

021

037

041

Fake news refers to deliberately spreading false or misleading information with the aim of deceiving the public, creating confusion, manipulating public opinion, or achieving specific political, economic, or social objectives. Online social networks (OSNs) have increased the convenience of realtime information dissemination, but they also lead to the rapid and widespread dissemination of fake news, causing detrimental effects on the online environment (Aïmeur et al., 2023). Detecting fake news has thus become a current research hotspot.

Early works primarily focused on manually extracting features from text content (Choudhary and Arora, 2021), such as the proportion of negation words, writing style, and language styles. However, traditional methods are inefficient and unable to handle large amounts of data. Therefore, researchers began to focus on deep learning-based automatic fake news detection. Bhattarai et al. (Bhattarai et al., 2021) captured the lexical and semantic properties of news text. Jin et al. (Jin et al.,



Figure 1: Motivation of our proposed MViR. We can see that different news texts describe the same image from various perspectives. For instance, some focus on the background building, others on the sign, and some on the person.

2016) detected fake news by leveraging significant disparities in image distributions.

043

047

049

051

054

060

061

062

063

064

065

066

067

069

070

With the development of OSNs (Aïmeur et al., 2023), multimodal fake news (Zhou and Zafarani, 2020; Singh et al., 2021), which includes text, images, and videos, has emerged. These forms are often more attractive and have a broader reach than traditional unimodal fake news. EANN (Wang et al., 2018) introduces an event discriminator to detect fake news. MVAE (Khattar et al., 2019) incorporates a multimodal variational autoencoder for multimodal fake news detection. MCAN (Wu et al., 2021) designs a co-attention network to better fuse multimodal features.

However, existing multimodal fake news detection methods struggle to effectively capture the multi-view visual-semantic relationships present in news content. News images often contain information from various perspectives. For instance, as shown in Figure 1, the left side illustrates an image associated with fake news, while the right side presents the accompanying text. It is evident that the text reports multiple aspects of the image from different viewpoints, a common phenomenon in real-world news articles. This poses a challenge for previous detection methods (Tufchi et al., 2023), as they often fail to consider the multi-view semantics of the content for trustworthy fake news detection.

We propose a multi-view visual-semantic rep-

1



Figure 2: The MViR framework consists of three modules: Multi-View Representation (MVR), Multi-View Feature Fusion (MVFF), and Multi-View Aggregation (MVA). It extracts image and text features, learns multi-view visual-semantic representations via MVR, fuses features with MVFF, and uses MVA to generate embeddings and predict fake news probabilities.

resentation for fake news detection (MViR) to address the above issues. Specifically, we propose a multi-view representation module to extract multiview fine-grained features from images, thereby providing the model with comprehensive multiview visual-semantic information. Afterward, we use a multi-view feature fusion module to fuse image and text information, further enhancing the representation capability of multi-view features. Finally, we use a Multi-View Aggregation module to process the fused features and extract multi-view semantic cues to enhance fake news detection. Our main contributions include:

071

072

084

880

100

101

102

104

106

108

(1) We design a multi-view representation module, which can explicitly model the multi-view semantics in news images and capture the multi-view features within the images.

(2) We design a multi-view aggregation module, which can explicitly learn multi-view embeddings, extract multi-view semantic cues from news and utilize them to enhance the model's ability to identify fake news.

(3) Experiments conducted on the widely used datasets show that MViR significantly outperforms previous approaches.

2 Methodology

As shown in Figure 2, MViR consists of three main parts: Multi-View Representation (MVR), Multi-View Feature Fusion (MVFF), and Multi-View Aggregation (MVA).

2.1 Feature Extraction

For a image I from multimodal news, we utilize VGG-19 (Simonyan and Zisserman, 2014) to extract visual features. These features are subsequently projected into a d-dimensional space via a fully connected (FC) layer, yielding image features represented as $V = [v_1, v_2, ..., v_r] \in \mathbb{R}^{r \times d}$, where r denotes the number of extracted regions. Similarly, for a text T containing m words, word embeddings are extracted using a pre-trained BERT (Devlin, 2018), which are mapped to the ddimensional space using a FC layer. Text features are expressed as $T = [t_1, t_2, ..., t_m] \in \mathbb{R}^{m \times d}$. 109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

2.2 Multi-View Representation

We propose a Multi-View Representation (MVR) module to capture the multi-view semantics of images. Understanding an image from different perspectives means that each region requires different attention from different perspectives. To achieve this, we employ a pyramid dilated convolution (Yu, 2015; Qu et al., 2020) layer with K parallel kernels to aggregate multi-scale contextual information from the image, which serves as the basis for calculating view-specific importance scores.

$$s_i^k = Convd(V, w^k, d^k), \quad k = 1, 2, \dots, K$$
 (1)

where w^k and d^k denote its kernel size and dilation rate, s_i^k denotes the output of the k-th kernel. We then concatenate these outputs:

$$s_i = Concat(s_i^1, \dots, s_i^K), \tag{2}$$

where $Concat(\cdot)$ denotes the concatenation of vectors. Afterward, a FC followed by a softmax activation is applied to compute the multi-view matrix $\hat{S} = [\hat{s}_1; ...; \hat{s}_r] \in \mathbb{R}^{r \times N}$, where N is the number of views, and \hat{s}_i is the *i*-th row vector. The above process can be summarized as follows:

$$\hat{s}_{ij} = \frac{\exp\left((W_s s_i + b_s)_j\right)}{\sum_{j=1}^r \exp\left((W_s s_i + b_s)_j\right)},$$
(3)

where $(W_s s_i + b_s) \in \mathbb{R}^{r \times N}$, $\hat{s}_i \in \mathbb{R}^N$ represents the importance scores of the *i*-th region over Nviews, $W_s \in \mathbb{R}^{N \times d}$ and $b_s \in \mathbb{R}^{1 \times N}$ are the learnable weights and bias, respectively. Finally, the image features can be summarized into a multiview representation $\mathbf{V}^* \in \mathbb{R}^{N \times d}$ as follows:

$$\mathbf{V}^* = \hat{\mathbf{S}}^T \mathbf{V}. \tag{4}$$

Dataset	Method	Accuracy	Fake News			Real News		
			Precision	Recall	F1 score	Precision	Recall	F1 score
	EANN (Wang et al., 2018)	0.827	0.847	0.812	0.829	0.807	0.843	0.825
	SAFE (Zhou et al., 2020)	0.816	0.818	0.815	0.817	0.816	0.818	0.817
	MCAN (Wu et al., 2021)	0.899	0.913	0.889	0.901	0.884	0.909	0.897
Waiha	CAFE (Chen et al., 2022)	0.840	0.855	0.830	0.842	0.825	0.851	0.837
weibo	FND-CLIP (Zhou et al., 2023)	0.907	0.914	0.901	0.907	0.917	0.901	0.908
	MSACA (Wang et al., 2024)	0.903	0.935	0.873	0.903	0.872	0.935	0.902
	EVENT-RADAR (Ma et al., 2024)	0.919	0.924	0.905	0.914	0.932	0.915	0.924
	MViR (Ours)	0.924	0.944	0.906	0.920	0.906	0.941	0.928
	EANN (Wang et al., 2018)	0.864	0.702	0.518	0.594	0.887	0.956	0.920
GossipCop	SAFE (Zhou et al., 2020)	0.838	0.758	0.558	0.643	0.857	0.937	0.895
	SPOTFAKE (Singhal et al., 2020)	0.858	0.732	0.372	0.494	0.866	0.962	0.914
	CAFE (Chen et al., 2022)	0.867	0.732	0.409	0.587	0.887	0.957	0.921
	FND-CLIP (Zhou et al., 2023)	0.880	0.761	0.549	0.638	0.899	0.959	0.928
	MSACA (Wang et al., 2024)	0.887	0.816	0.538	0.648	0.897	0.971	0.933
	RaCMC (Yu et al., 2024)	0.879	0.745	0.563	0.641	0.902	0.954	0.927
	MViR (Ours)	0.895	0.784	0.619	0.692	0.914	0.963	0.937

Table 1: Results on two datasets. The best performance is in bold, while underlining highlights the follow-up.

144 145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

168

169

170

171

172

173

174

175 176

2.3 Multi-View Feature Fusion

We propose a Multi-View Feature Fusion (MVFF) module that combines multi-view image features with the corresponding text features, further enhancing their representational capacity. MVFF consists of l multi-view fusion layers, each containing a co-attention mechanism (Lu et al., 2019) and a feed-forward network (FFN). Both components are enclosed by a residual connection and followed by layer normalization. The co-attention extends the standard multi-head attention by using queries (Q) from one modality and keys (K) and values (V) from another. In our approach, Q comes from the multi-view features of the image or fused features, while K and V come from the text.

$$Q_i = \mathbf{V}^* W_i^Q, \quad K_i = T W_i^K, \quad V_i = T W_i^V, \quad (5)$$

where $W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{1 \times d_h}$ are the projection matrices for the *i*-th head, H denotes the number of heads, $d_h = d/H$ is the dimension of the output feature of each head. The calculation process of the co-attention can be presented as follows:

$$MultiHead(Q, K, V) = Concat(h_1, h_2, ..., h_H)W^O + V$$
(6)

where $W^O \in \mathbb{R}^{d \times d}$ is learnable weights, and $h_i = Att(Q_i, K_i, V_i)$. Att denotes the scaled-dot product attention, defined as follows:

$$Att(Q_i, K_i, V_i) = softmax\left(\frac{Q_i K_i^{\top}}{\sqrt{d_h}}\right) V_i.$$
(7)

To enhance the representational capacity of the fused features, the output of the co-attention is processed through an FFN. It is implemented as a twolayer multi-layer perceptron (MLP) with a ReLU activation function applied between the layers. The above process can be summarized as follows:

$$\mathbf{X} = FFN(MultiHead(Q, K, V)) \oplus Q, \qquad (8)$$

177where $\mathbf{X} \in \mathbb{R}^{N \times d}$ denotes the features after fu-178sion, \oplus represents the fusion operation, e.g., vector179concatenation or elementwise add.

2.4 Multi-View Aggregation

The uniqueness of the Multi-View Aggregation (MVA) module lies in its use of multiple aggregators to generate a set of embeddings from the fused features, explicitly modeling multi-view features. This allows for the evaluation of news authenticity through multi-view semantic cues. Specifically, after obtaining a fused feature set $\{\mathbf{x}_n\}_{n=1}^N$, a series of feature aggregators $\{f_n^v\}_{n=1}^N$ are used to aggregate $\{\mathbf{x}_n\}_{n=1}^N$ into a set of semantic cues embeddings $\{\hat{\mathbf{x}}_n\}_{n=1}^N$:

$$\hat{\mathbf{x}}_n = f_n^v \left(\{ \mathbf{x}_n \}_{n=1}^N \right), \tag{9}$$

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

202

204

205

207

208

209

210

211

212

213

214

215

216

where $\hat{\mathbf{x}}_n \in \mathbb{R}^d$. Each $\hat{\mathbf{x}}_n$ represents the semantic cues of a set of views. N is the number of views.

Next, we combine the semantic cues from each view with the text features and use a decision network to assess whether the news is true or fake. We employ only a single aggregator to obtain the text embedding:

$$z_n = \max\left(0, W_f \operatorname{Concat}(\hat{\mathbf{x}}_n, f^t(T))\right), \qquad (10)$$

$$\hat{\boldsymbol{y}} = \max_{n=1}^{N} \left(\operatorname{softmax} \left(\boldsymbol{z}_n \boldsymbol{W}_l \right) \right), \tag{11}$$

where \hat{y} represents the probability of the news being fake, W_f represents the parameters of the fully connected layers, and W_l represents the parameters of the linear layer within the softmax function. For multiple semantic cues from different views of news, if any of these cues is detected as fake, the news is classified as fake.

2.5 Objective Function

We leverage cross-entropy to measure the classification loss and train our model:

$$\mathcal{L} = \sum_{i=1}^{\mathcal{N}} - \left[y_i * \log\left(\hat{y}_i\right) + (1 - y_i) * \log\left(1 - \hat{y}_i\right) \right], \quad (12)$$

where \mathcal{N} denotes the number of news reports, y_i represents the ground-truth label of the *i*-th news. Labels 0 and 1 refer to real news and fake news, respectively.

Table 2: Ablation study on two datasets.						
Method		Accuracy	F1 score			
		Accuracy	Fake News	Real News		
	MViR (Ours)	0.924	0.920	0.928		
	w/o MVR	0.901	0.893	0.909		
Waiha	w/o MVFF	0.894	0.884	0.902		
weibo	w/o MVA	0.907	0.904	0.909		
	Max Probability(Real)	0.912	0.908	0.915		
	Average Probability	0.918	0.915	0.921		
	MViR (Ours)	0.895	0.692	0.937		
	w/o MVR	0.883	0.648	0.930		
CassinCan	w/o MVFF	0.881	0.667	0.927		
GossipCop	w/o MVA	0.886	0.639	0.932		
	Max Probability(Real)) 0.874	0.655	0.915		
	Average Probability	0.884	0.653	0.929		

Table 3: Analysis for different numbers of MVFF layer.

	Lovers	Accuracy	F1 score		
	Layers	Accuracy	Fake News	Real News	
	2	0.912	0.908	0.915	
Weibo	3	0.924	0.920	0.928	
	4	0.917	0.915	0.921	
	2	0.884	0.648	0.930	
GossipCop	3	0.895	0.692	0.937	
	4	0.891	0.674	0.935	

3 Experiments

3.1 Datasets and Experimental Settings

We performed experiments using two real-world datasets: Weibo (Jin et al., 2017) and GossipCop (Shu et al., 2020). The Weibo dataset contains a total of 9,528 news, with 7,532 used for training and 1,996 used for testing. The GossipCop dataset contains a total of 12,840 news, with 10,010 used for training and 2,830 used for testing. For fairness in comparison, we adhered to the data-splitting protocol and processing steps used in prior works (Wu et al., 2021; Ying et al., 2023).

Experiments were carried out on an NVIDIA Tesla A100 GPU. For the Weibo and GossipCop datasets, we adopted bert-base-chinese and bertbase-uncased respectively to extract text features. More details of the implementation can be found in the appendix A.

3.2 Performance Comparison

Table 1 shows the performance comparison between MViR and the baseline methods. MViR demonstrated excellent performance across all metrics, including accuracy, precision, recall, and F1 score. MViR achieved an average accuracy of 92.4% on the Weibo dataset and 89.5% on the GossipCop dataset, outperforming the best existing models by 1.9 and 1.7 percentage points.

While many methods, such as MCAN, detect fake news by fusing multimodal features, they do not consider the multi-view characterization of fake news. In contrast, MViR effectively captures the multi-view features of images, thus improving the performance of multimodal fake news detection.



250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

268

269

270

271

272

273

274

275

276

277

278

279

281

282

284

285

286

288

289

290

291

292

293

294

3.3 Ablation Study

To evaluate the contribution of each component in MViR, we removed the MVR, MVFF, and MVA modules individually. Table 2 shows that removing any module leads to a significant performance drop, confirming that the components complement each other to improve fake news detection. Additionally, we compared decision networks using maximum real news probability, average probability, and maximum fake news probability. The results indicate that the network using maximum fake news probability performs best, as it can capture more reliable key features of fake news.

3.4 Parameter Sensitivity Analysis

Impact of the Number of Feature Fusion Layers: As shown in Table 3. The experimental results demonstrate that appropriately increasing the number of layers can promote interactions between different modality features, thereby enhancing model performance. However, when the number of layers becomes larger, further increases result in performance degradation, likely due to excessive model complexity leading to overfitting.

Impact of the Number of views: We also analyzed the impact of the number of viewpoints on model performance, as shown in Figure 3. It can be observed that appropriately increasing the number of viewpoints helps capture richer details and more diverse features, thereby improving performance. However, when the number of viewpoints increases further, the model's performance may decline due to the introduction of redundant information, which could affect its generalization ability. In our experiments, MViR performed best with 12 viewpoints.

4 Conclusion

In this work, we propose MViR, a novel framework for fake news detection using multi-view visual-semantic representations. Our approach includes a multi-view representation module to extract visual-semantic features from images, a feature fusion module to combine image and text features, and a Multi-View Aggregation module to learn multi-view embeddings. Experiments on two benchmark datasets show that MViR outperforms existing state-of-the-art methods.

241

242

243

245

246

247

249

217

218

295 Limitations

301

304

307

308

312

313

314

315

317

319

321

322

323

337

338

340

341

In this work, we propose MViR for detecting multimodal fake news. Although our approach demonstrates excellent performance, it still has the following limitations:

- Although MVIR performs best on most metrics, it underperforms EVENT-RADAR and MSACA on a few individual metrics (e.g., real news precision on the Weibo dataset). This may stem from dataset-specific biases or the model's sensitivity to certain semantic cues, which warrants further investigation.
 - Experiments were conducted on the Weibo and GossipCop datasets, which are relatively small and domain-specific. The model's generalization capability on larger, cross-domain datasets (e.g., Twitter, news articles) requires further validation.
 - The multimodal processing pipeline introduces additional complexity, but the computational overhead (e.g., training time, memory usage) has not been quantified. Efficiency analysis is critical for practical deployment, especially on social platforms requiring realtime detection. In future work, we will further evaluate MViR's efficiency and conduct more comprehensive experiments on additional datasets.

Ethical Statement

Our proposed multimodal detection method aims to reduce false news on the internet and maintain a healthy online environment. This paper adheres 326 to the ACM Code of Ethics and Professional Con-327 duct. All experimental data used in our research originates from publicly accessible resources, and 329 before utilization, we confirm that it complies with relevant usage regulations and does not contain 331 sensitive private information. Additionally, appropriate citations are given to the sources of related 333 papers and pre-trained models utilized in our work. 334 Lastly, our code will be released under the license applicable to any artifacts used. 336

References

- Esma Aïmeur, Sabrine Amri, and Gilles Brassard. 2023. Fake news, disinformation and misinformation in social media: a review. *Social Network Analysis and Mining*, 13(1):30.
- Liesbeth Allein, Marie-Francine Moens, and Domenico 342 Perrotta. 2021. Like article, like audience: Enforcing 343 multimodal correlations for disinformation detection. 344 arXiv preprint arXiv:2108.13892. 345 Bimal Bhattarai, Ole-Christoffer Granmo, and Lei Jiao. 346 2021. Explainable tsetlin machine framework for 347 fake news detection with credibility score assessment. 348 arXiv preprint arXiv:2105.09114. 349 Yixuan Chen, Dongsheng Li, Peng Zhang, Jie Sui, Qin 350 Lv, Lu Tun, and Li Shang. 2022. Cross-modal am-351 biguity learning for multimodal fake news detection. 352 In WWW. 353 Anshika Choudhary and Anuja Arora. 2021. Linguistic 354 feature based learning model for fake news detection 356 and classification. Expert Systems with Applications, 169:114171. 357 Jacob Devlin. 2018. Bert: Pre-training of deep bidi-358 rectional transformers for language understanding. 359 arXiv preprint arXiv:1810.04805. 360 Zhiwei Jin, Juan Cao, Han Guo, Yongdong Zhang, and 361 Jiebo Luo. 2017. Multimodal fusion with recurrent 362 neural networks for rumor detection on microblogs. 363 In ACM MM. 364 Zhiwei Jin, Juan Cao, Yongdong Zhang, Jianshe Zhou, 365 and Qi Tian. 2016. Novel visual and statistical im-366 age features for microblogs news verification. IEEE 367 TMM. 368 Dhruv Khattar, Jaipal Singh Goud, Manish Gupta, and 369 Vasudeva Varma. 2019. Mvae: Multimodal varia-370 tional autoencoder for fake news detection. In WWW. 371 Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 372 2019. Vilbert: Pretraining task-agnostic visiolin-373 guistic representations for vision-and-language tasks. 374 NeurIPS. 375 Zihan Ma, Minnan Luo, Hao Guo, Zhi Zeng, Yiran Hao, 376 and Xiang Zhao. 2024. Event-radar: Event-driven 377 multi-view learning for multimodal fake news detec-378 tion. In Proceedings of the 62nd Annual Meeting of 379 the Association for Computational Linguistics (Volume 1: Long Papers), pages 5809-5821. 381 Leigang Qu, Meng Liu, Da Cao, Liqiang Nie, and 382 Qi Tian. 2020. Context-aware multi-view summa-383 rization network for image-text matching. In Pro-384 ceedings of the 28th ACM international conference on multimedia, pages 1047-1055. 386 Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dong-387 won Lee, and Huan Liu. 2020. Fakenewsnet: A data 388 repository with news content, social context, and spa-389 tiotemporal information for studying fake news on 390 social media. Big data, 8(3):171-188. 391 Karen Simonyan and Andrew Zisserman. 2014. Very 392 deep convolutional networks for large-scale image 393 recognition. arXiv preprint arXiv:1409.1556. 394

Vivek K Singh, Isha Ghosh, and Darshan Sonagara. 2021. Detecting fake news stories via multimodal analysis. *Journal of the Association for Information Science and Technology*, 72(1):3–17.

398

400 401

402

403

404

405

406

407 408

409

410

411

412 413

414

415

416

417 418

419 420

421

422

423

424

425

426

427

428

429

430

431 432

433 434

435

436

437

438

439 440

441

442

443

- 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 *AAAI*.
- Shivani Tufchi, Ashima Yadav, and Tanveer Ahmed. 2023. A comprehensive survey of multimodal fake news detection techniques: advances, challenges, and opportunities. *International Journal of Multimedia Information Retrieval*, 12(2):28.
- Jiandong Wang, Hongguang Zhang, Chun Liu, and Xiongjun Yang. 2024. Fake news detection via multiscale semantic alignment and cross-modal attention. In *Proceedings of the 47th International ACM SI-GIR Conference on Research and Development in Information Retrieval*, pages 2406–2410.
 - 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 *KDD*.
 - Yang Wu, Pengwei Zhan, Yunjian Zhang, Liming Wang, and Zhen Xu. 2021. Multimodal fusion with coattention networks for fake news detection. In *ACL-IJCNLP*.
 - Qichao Ying, Xiaoxiao Hu, Yangming Zhou, Zhenxing Qian, Dan Zeng, and Shiming Ge. 2023. Bootstrapping multi-view representations for fake news detection. In *AAAI*.
 - F Yu. 2015. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*.
- Xinquan Yu, Ziqi Sheng, Wei Lu, Xiangyang Luo, and Jiantao Zhou. 2024. Racmc: Residual-aware compensation network with multi-granularity constraints for fake news detection. *arXiv preprint arXiv:2412.18254*.
- Xinyi Zhou, Jindi Wu, and Reza Zafarani. 2020. : Similarity-aware multi-modal fake news detection. In *Pacific-Asia Conference on knowledge discovery and data mining*, pages 354–367. Springer.
- Xinyi Zhou and Reza Zafarani. 2020. A survey of fake news: Fundamental theories, detection methods, and opportunities. *ACM Computing Surveys*, 53(5):1–40.
- Yangming Zhou, Yuzhou Yang, Qichao Ying, Zhenxing Qian, and Xinpeng Zhang. 2023. Multimodal fake news detection via clip-guided learning. In *ICME*.

445 446 447 448 449 450 451 452 453 454 455 456 457 458 459 460 461 462 463

464

465

466

467

468

469

470

471

472

473

474

475

444

A Implementation

We implemented MViR using PyTorch 2.3.1 and conducted all experiments on a single NVIDIA Tesla A100 GPU. For text feature extraction, we used bert-base-chinese for the Weibo dataset with a maximum sequence length of 160, and bert-baseuncased for the GossipCop dataset with a maximum sequence length of 394. Images were resized to 224×224 to match the input dimensions of the pre-trained VGG-19 model. The dimensions of image and text features *d* were set to 256, with the number of heads set to 4 and a dropout rate of 0.5. We trained the model using AdaBelief for 50 epochs with a batch size of 32 and an initial learning rate of 1e-4. Additional implementation details can be found in the code.

The specific Pyramid Dilated Convolutional Layers are shown in Table 4, where k denotes the kernel number, w^k represents the kernel size, d^k represents the dilation rate, and s^k represents output channel of k-th convolution kernel, respectively. As the dilation rate increases, the receptive field of the kernel is enlarged without reducing the regional resolution.

We also provide a statistical overview of the detailed parameters for the two datasets (Weibo and GossipCop) used in our study, as shown in Table 5.

Table 4: Configurations of pyramid dilated convolution

k	1	2	3	4	5	6	7
w^k	1	3	3	3	5	5	5
d^k	1	1	2	3	1	2	3
s^k	256	128	128	128	128	128	128

	Weibo	GossipCop
Total news	9528	12840
Images	13272	15488
Fake news	3783(Train)	2036(Train)
	1000(Test)	545(Test)
Real news	3749(Train)	7974(Train)
	996(Test)	2285(Test)

B Baselines

To evaluate the performance of MViR, we compared it with several baseline approaches in our experiments, all of which are classic schemes in the field of fake news detection. We provide a brief introduction to each of them:

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

EANN (Wang et al., 2018) enhances the detection capability of fake news in new events by building an end-to-end framework that includes multimodal feature extraction, fake news detection, and event discrimination to learn event-invariant feature representations.

SAFE (Zhou et al., 2020) effectively identifies fake news by jointly learning the textual and visual features of news articles and their cross-modal relationships, utilizing a similarity-aware approach to detect content mismatches.

MCAN (Wu et al., 2021) better fuses textual and visual features for fake news detection by learning the inter-dependencies among multimodal features.

CAFE (Chen et al., 2022) is a method for detecting fake news that improves accuracy by assessing ambiguities between different media types and capturing how they relate to each other.

FND-CLIP (**Zhou et al., 2023**) is a framework that uses CLIP technology to combine text and image information for better fake news detection.

MSACA (Wang et al., 2024) is a network that improves fake news detection by aligning text and images at multiple scales and using attention to select the best features.

EVENT-RADAR (Ma et al., 2024) is a framework that detects fake news by analyzing multimodal information in events and calculating the credibility of each view.

SPOTFAKE (Singhal et al., 2020) is a method that uses transfer learning to combine text and image information for more accurate fake news detection.

RaCMC (Allein et al., 2021) enhances the differences between real and fake news by utilizing multi-scale feature interaction and fusion, along with multi-granularity constraints, to improve the accuracy of fake news detection.

C Case Study

To further illustrate the importance of multiperspective analysis for fake news detection, we compared the detection results of MViR with baseline approache MVAE and showcased some fake news instances (translated into English) that were correctly identified by MViR but overlooked by MVAE, as shown in the Figure 4. A common feature of these fake news items is their rich image perspective information; previous works, which ig-



Figure 4: Case study.

nored the multi-perspective characteristics of fake 526 news, resulted in inaccurate identification. In con-527 trast, MViR successfully recognized these fake 528 news instances by capturing multi-perspective fea-529 tures from both images and text. These exam-530 531 ples demonstrate that a single image representation struggles to comprehensively describe image infor-532 mation, leading to misclassification of certain fake 533 news items that contain multi-perspective details. 534