Structure-adaptive Adversarial Contrastive Learning for Multi-Domain Fake News Detection

Anonymous ACL submission

Abstract

The rapid proliferation of fake news across multiple domains poses significant threats to society. Existing multi-domain detection models typically capture domain-shared semantic fea-004 tures to achieve generalized detection. However, they often fail to generalize well due 007 to poor adaptability, which limits their ability to provide complementary features for detection, especially in data-constrained conditions. To address these challenges, we investigate the propagation-adaptive multi-domain fake news detection paradigm. We propose a 012 novel framework, Structure-adaptive Adversarial Contrastive Learning (StruACL), to adap-015 tively enable structure knowledge transfer between multiple domains. Specifically, we first 017 contrast representations between content-only and propagation-rich data to preserve structural patterns in the shared representation space. Additionally, we design a propagation-guided adversarial training strategy to enhance the diversity of representations. Under the StruACL objective, we leverage a unified Transformerbased and graph-based model to jointly learn transferable semantic and structural features for detection across multiple domains. Experiments on seven fake news datasets demon-027 strate that StruACL-TGN achieves better multidomain detection performance on general and data-constrained scenarios, showing the effectiveness and better generalization of StruACL.

1 Introduction

033

037

041

Nowadays, mainstream social platforms have facilitated the news dissemination in a faster and cheaper way. Nevertheless, the ease has also caused the wide spread of fake news, which has had detrimental effects on individuals and society (Loomba et al., 2021). Triggered by the negative impact of fake news, fake news detection has become a pressing challenge due to its widespread impact across diverse platforms and domains.

News content and its corresponding user engagements (i.e., tree-structured propagation) are two key data types in detecting fake news. Contentbased detection methods (Ma et al., 2016; Ruchansky et al., 2017; Karimi and Tang, 2019) capture intrinsic semantic or linguistic features of claim tweets to detect fake news. Propagation-based detection methods (Ma et al., 2018; Kumar and Carley, 2019; Ma and Gao, 2020; Hu et al., 2021; Bian et al., 2020; Wei et al., 2021; Lin et al., 2021; Wei et al., 2022a) are designed to integrate structural features to complement textual content for detection. Nevertheless, in real-world scenarios, labeled data for fake news is often scarce, particularly in specific low-resource domains or emerging topics, which hinders detection performance. Recently, multi-domain fake news detection has been widely studied to leverage and integrate knowledge from multi-domain data to improve target-domain detection (Zhu et al., 2023; Liang et al., 2022; Wang et al., 2018; Zhang et al., 2021; Nan et al., 2021; Li et al., 2024), alleviating the data limitation challenge to some extent.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

However, the representations learned by most existing multi-domain detection paradigms fail to generalize well due to poor adaptability to the propagation structure. Firstly, as shown in Figure 1, some multi-domain approaches primarily focus on learning domain-invariant or domain-shared semantic features (Wang et al., 2018; Nan et al., 2021) on content-only training data. However, semantic features inherently differ from structural patterns, rendering these content-based methods inadequate for generalizing to samples that involve propagation. Furthermore, directly extending domainspecific propagation-based methods struggles to effectively adapt to detection scenarios lacking propagation structures, resulting in suboptimal detection performance for content-only samples (Wei et al., 2024). Therefore, a critical challenge lies in learning more robust representations by enhancing



Figure 1: Difference between *propagation-adaptive multi-domain fake news detection* in this study and existing multi-domain fake news detection paradigms.

structure adaptability for multi-domain fake news detection.

In this paper, we study a novel propagationadaptive multi-domain fake news detection paradigm, where the detection model is trained on both propagation-based data and content-only data. Our goal is to enhance generalization for both types of input.

To achieve this, we propose a new propagation structure-adaptive adversarial contrastive learning framework (StruACL) to adaptively learn generalized semantic and structural representations for multi-domain fake news detection. Specifically, we first design a new structure-aware contrastive learning (StruCL) objective to facilitate the adaptive transfer of structure knowledge during multidomain training. With the guidance of structure label, StruCL leverages contrastive learning to differentiate representations between samples with and without propagation, effectively capturing and retaining structural knowledge in the shared representation space. By integrating this structural information, the learned representations become more informative, allowing the model to achieve enhanced performance in detecting fake news across both propagation-based and content-only domains. Additionally, we design a propagation-guided adversarial training (PAT) strategy to enhance the diversity of representations under the data-constrained condition. PAT adaptively performs adversarial perturbations on original embeddings using the Fast Gradient Method (FGM) (Miyato et al., 2017) to generate worst-case samples for both content-only and propagation-based inputs. By jointly contrasting on both original and adversarial samples, the model can further effectively learn fine-grained semantic and structure knowledge via retaining the propagation-adaptive feature consistency. For the model architecture, we adopt a shared Transformerbased and graph-based network to jointly encode semantic and structural features from news content and available propagation across multiple domains, respectively. Under the proposed objective, our StruACL-TGN generalizes well across both content-only and propagation-structured domains. 122

123

124

125

126

127

128

129

130

131

132

133

134

135

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

159

We conduct experiments on seven fake news datasets with and without propagation. The experimental results demonstrate that our StruACL-TGN achieves superior performance in multi-domain fake news detection. Extensive experiments show the effectiveness of StruACL objective, particularly in data-limited application scenarios.

The main contributions are as follows: 1) We study a novel propagation-adaptive multi-domain fake news detection paradigm and develop a novel StruACL-TGN to learn generalized representations for detection on both domains with propagation data and content-only domains. 2) We design a new StruACL framework to learn more informative multi-domain representations. It contrasts semantic and structural representations to preserve and transfer structural knowledge, as well as introduces propagation-guided adversarial training to enhance the diversity of representations. 3) Experiments on seven fake news datasets demonstrate that StruACL-TGN achieves superior multi-domain detection performance. Extensive experiments further show that StruACL enhances the model's generalization capabilities in data-constraint applications.

2 Methodology

In this section, we first describe the problem definition of propagation-adaptive multi-domain fake news detection. Then, we propose a new StruACL-TGN to learn generalized representations on both domains with propagation data and content-only domains. The overall architecture is shown in Figure 2. It adopts a shared Transformer-based

119

121



Figure 2: Overview of the proposed StruACL-TGN.

and graph-based network to encode semantic and 160 structural knowledge across multiple domains. For 161 model training, we first propose a new StruCL to effectively utilize structure knowledge. Additionally, 163 we design propagation-guide adversarial training 164 (PAT) that generates worst-case samples to enhance 165 the diversity of representations. Through applying PAT on original and adversarial samples, our Stru-167 ACL can learn more informative multi-domain representations from domains with propagation data 169 and domains with content-only data. 170

2.1 Problem Definition

171

172

173

174

176

177

178

181

183

184

188

190

Unlike existing multi-domain fake news detection tasks, propagation-adaptive multi-domain fake news detection aims to detect fake news across domains with heterogeneous data availability, where some domains have both propagation structures and content, while others only have content.

Formally, let \mathcal{K} represent the set of all domins. Define $D^{(k)}$ as the dataset of each domain $k \in \mathcal{K}$. For each domain k that includes propagation data, $D^{(k)}$ is defined as $D^{(k)} = \{(x_i^{(k)}, G_i^{(k)}, y_i^{(k)})\}_{i=1}^{N_k}$, where $x_i^{(k)}$ is the content of the *i*-th news sample in domain k. $G_i^{(k)}$ is the propagation structure (e.g., a tree-like graph) associated with $x_i^{(k)}$. $y_i^{(k)} \in \{0, 1\}$ is the label indicating whether the news is fake or real. N_k is the number of samples in domain k. For each domain k' that excludes propagation data, $D^{(k)}$ is defined as $D^{(k')} = \{(x_i^{(k')}, y_i^{(k')})\}_{i=1}^{N_{k'}}$. where $x_i^{(k')}$ is the content of the *i*-th news sample in domain k'. $y_i^{(k')} \in \{0, 1\}$ is the label. $N_{k'}$ is the number of samples in domain k'. Propagation-adaptive multi-dom

Propagation-adaptive multi-domain fake news detection aims to utilize both rich propagation structures and content-only samples to enhance detection performance across various domains. Fake news detection can be regarded as a binary classification task. Specifically, the objective is to learn a unified detection model $f(\cdot)$ that predicts the label \hat{y} (e.g., fake or real) for a news item x (with or without propagation G across all domains: $\hat{y} = f(x, G)$, where $G = \emptyset$ for domains without propagation data.

191

192

193

194

195

200

201

202

203

204

206

207

208

209

211

212

213

214

215

216

217

218

219

221

2.2 Model Architecture

The network structure consists of a shared Transformer-based semantic encoder, a graphbased structure encoder, and a hybrid fake news classifier.

Transformer-based Semantic Encoder Considering multilingual settings, a pretrained multilingual BERT model (Conneau et al., 2020) on a monolingual corpus is utilized to facilitate language adaptation. Formally, given an input token sequence $x_{i1}, ..., x_{iT}$ where x_{ij} refers to *j*-th token in the *i*-th input sample, and *T* is the maximum sequence length, the model learns to generate the context representation of the input token sequences:

$$\mathbf{h}_{i}^{s} = \text{BERT}([\text{CLS}], x_{i1}, ..., x_{iN}, [\text{SEP}]), (1)$$

where [CLS] and [SEP] are special tokens, typically placed at the beginning and end of each sequence, respectively. \mathbf{h}_i^s indicates the hidden representation of the *i*-th input sample, computed by the representation of [CLS] token in the last layer of the encoder.

225

227

237

240

241

242

243

244

245

247

249

251

258

260

261

262

264

Graph-based Structure Encoder Based on the semantic representations, propagation-based models integrate structural features to enhance detection. Graph neural networks are widely applied to extract structural features through message-passing across nodes in the propagation graph. Given the input sample, which includes the textual content of the source news x and propagation trees G, existing models utilize various neural networks to extract high-level textual and structural features for detection. The formulation is defined as,

$$\mathbf{h}_{i}^{g} = \mathrm{GNN}(x_{i}, G_{i}; \Theta), \tag{2}$$

where $GNN(\cdot)$ refers to graph-based encoders in propagation-based models (Bian et al., 2020; Wei et al., 2022b), and Θ refers to the corresponding trainable parameters. The input embedding of x_i and context c_i in G are initialized with the semantic embedding \mathbf{h}_i^s .

Hybrid Fake News Classifier To address the feature gap between different domains in distinguishing fake news, we design a hybrid fake news classifier to learn domain-specific and domain-shared discriminative features for detection. Specifically, based on the final representation, domain-specific fake news classifiers are employed to predict the veracity label of each news content. For domain k, the initial prediction distribution is computed as,

$$\hat{\mathbf{y}'}^{(k)} = \mathbf{W}_c^{(k)} \mathbf{z} + \mathbf{b}_c^{(k)}, \qquad (3)$$

where $\mathbf{W}_{c}^{(k)}$ and $\mathbf{b}_{c}^{(k)}$ are trainable parameters of domain k's classifier. Similarly, we apply a parameter-shared classifier to predict the veracity label for all domains, i.e.,

$$\hat{\mathbf{y}}^s = \mathbf{W}_c \mathbf{z} + \mathbf{b}_c. \tag{4}$$

Based on the above prediction, the final prediction for domain k is defined as,

$$\hat{\mathbf{y}}^{(k)} = Softmax(\frac{\hat{\mathbf{y}}^s}{2} + \frac{\hat{\mathbf{y}'}^{(k)}}{2}).$$
(5)

2.3 Optimization Process

2.3.1 Classification Objective

To achieve fake news detection, the model is trained by minimizing a joint loss function across all domains, considering both content and propagation data, i.e.,

$$\mathcal{L}_{\text{CLS}} = \sum_{k \in \mathcal{K}_P} \frac{1}{N_k} \sum_{i=1}^{N_k} \ell(f(x_i^{(k)}, G_i^{(k)}), y_i^{(k)}) + \sum_{k' \in \mathcal{K}_C} \frac{1}{N_{k'}} \sum_{i=1}^{N_{k'}} \ell(f(x_i^{(k')}, \emptyset), y_i^{(k')})$$

$$(2)$$

where $\mathcal{K}_P \subseteq \mathcal{K}$ is the set of domains with propagation data, and $\mathcal{K}_C \subseteq \mathcal{K}$ is the set of domains without propagation data. $\ell(\cdot, \cdot)$ is the cross-entropy classification loss.

2.3.2 Structure-aware Contrastive Learning

We design a new structure-aware contrastive learning (StruCL) objective to facilitate the adaptive transfer of structure knowledge during multidomain training. It contrasts between representations with and without propagation structures. The objective of StruCL is defined as,

$$\begin{aligned} \mathcal{L}_{\text{StruCL}} &= \sum_{k \in \mathcal{K}_{P}} \mathcal{L}_{\text{StruCL}}^{(k)} \\ &= \sum_{k \in \mathcal{K}_{P}} -\frac{1}{N_{k}} \sum_{i=1}^{N_{k}} \left(\log \frac{e^{sim(\mathbf{z}_{i}^{(k)}, \mathbf{z}_{pos}^{(k)})/\tau}}{e^{sim(\mathbf{z}_{i}^{(k)}, \mathbf{z}_{pos}^{(k)})/\tau} + \sum_{j=1}^{N_{k}} \mathcal{K}_{g_{i}^{(k)} \neq g_{j}^{(k)}} e^{sim(\mathbf{z}_{i}^{(k)}, \mathbf{z}_{j}^{(k)})/\tau}} \right), \end{aligned}$$
(6)

where $g_i = 1$ refers to the hidden representation with structure information. The indicator function $\mathbb{F}_{q_i \neq q_i}$ equals 1 when the propagation structure label g_i and g_j are different, indicating a negative sample. $sim(\cdot, \cdot)$ is a pairwise similarity function, i.e., dot product. $\tau > 0$ is a scalar temperature parameter that controls the separation between the class with and without propagation structure. By minimizing the objective loss, the model is encouraged to separate semantic and structural representations while bringing closer the representations of the same type of input. This approach preserves structural features in the shared representation space, enhancing the model's ability to generalize and detect samples of different input types during testing.

2.3.3 Propagation-guided Adversarial Training

At each step of training, under structure-aware contrastive learning and multi-domain classification objectives, we apply an adversarial training strategy (e.g., FGM (Miyato et al., 2017)) on original samples to produce adversarial perturbations. Specifically, the perturbations are put on the embedding layers of semantic encoder, and then obtain adversarial samples. After that, we leverage the joint 265

268 269

270

271

272

273

274

275

276

278

279

281

282

283

284

289

290

291

292

293

294

295

296

297

299

300

301

302

303

304

390

391

392

393

394

395

396

398

objective on these worst-case samples to maximize 305 the consistency of transferable representations with 306 or without propagation across multiple domains. 307 Under the joint objective on both original and adversarial samples, our model can learn propagationrobust transferable features for multi-domain fake news detection. The optimization objective for cor-311 responding adversarial samples can be derived by 312 following the calculation process for the original 313 samples, denoted as, $\mathcal{L}_{CLS}^{r-adv} + \mathcal{L}_{StruCL}^{r-adv}$. Take the 314 domain with propagation data as an example, the 315 adversarial perturbation for content-only samples 316 is defined as, 317

$$\min_{\theta} \mathbb{E}_{(x^{(k)}, y^{(k)}) \sim D^{(k)}} \max_{\|r_{\text{adv}}\|_q \le \epsilon} \left(L_{\text{CLS}} + L_{\text{StruCL}} \right)$$
(7)

(7) where $r_{adv} = -\epsilon \frac{g}{\|g\|_q}, g = \nabla \log p(y^{(k)}|x^{(k)}; \hat{\theta}).$ The overall loss of StruACL is defined as a sum

of joint objective on both original and adversarial samples, i.e.,

 $\mathcal{L}_{total} = \mathcal{L}_{CLS} + \mathcal{L}_{StruCL} + \mathcal{L}_{CLS}^{r-adv} + \mathcal{L}_{StruCL}^{r-adv}.$ (8)

3 Experimental Setups

3.1 Datasets

318

319

321

322

324

325

326

328

332

333

334

336

338

339

341

345

347

348

351

We conduct experiments on seven widely-used public datasets for fake news detection, where two content-based fake news datasets including Weibo21 (Nan et al., 2021), and Covid19 (Patwa et al., 2021), and five propagation-based fake news datasets including Twitter, TwitterCovid19 (Kar et al., 2021; Lin et al., 2022), WeiboCovid19 (Lin et al., 2022), Arabic (Alam et al., 2021; Lin et al., 2023), and Cantonese (Ke et al., 2020; Lin et al., 2023). Based on the above seven datasets, we build two major benchmarks to achieve different multi-domain fake news detection settings, each involving at least one content-based and propagationbased datasets to evaluate potential transferability of detection methods between semantic and propagation structure. Specifically, CovidEval includes five datasets related to the same event (i.e., COVID-19): Covid19, WeiboCovid19, Twitter-Covid19, Arabic, and Cantonese. CrossEval includes three datasets from different social platforms and social events: Weibo21, Twitter, and Weibo-Covid19. Please see appendix for detailed description and statistics.

3.2 Evaluation Metrics

Since fake news detection can be regarded as a binary classification, we adopt widely-used evalua-

tion metrics for classification task, including accuracy (ACC), macro-averaged F1 score (F1).

3.3 Comparison Methods

We compare with single-domain fake news detection methods, and four multi-domain detection methods. *Single-domain methods*. **XLM-RoBERTa** (Conneau et al., 2020) uses a PLMbased semantic encoder with a linear classifier for fake news detection. **TextCNN** (Kim, 2014) utilizes convolutional layers to extract local semantics from news content.

Multi-domain methods. XLM-RoBERTa-M is an extended version of XLM-RoBERTa designed for multi-domain detection via parameter sharing across multiple domains. EANN (Wang et al., 2018) learns domain-invariant representations for detection. We re-implement by only considering the textual modality of news content across multiple domains. **MDFEND** Nan et al. (2021) uses a domain gate to aggregate multi-domain semantic representations. SAT (Wei et al., 2024) learns structure-invariant features from samples with and without propagation for detection. We extend this framework for multi-domain detection with the same model architecture of our method, denoted as SAT-TGN. More details of the related works are listed in the Appendix.

3.4 Implementation Details

All experiments are conducted on a single NVIDIA Tesla A100 80GB card. We use multilingual pretrained models to extract textual features considering different languages across datasets, and finetune the semantic encoder during training. The dimension of hidden vectors is set to 64. The graph layers are set to 2. The learning rate is set to 0.0001. The Adamax optimizer is adopted for all methods with the learning rate initialized to 0.0001 and weight decay as 0. The temperature parameter is searched from $\{0.1, 1\}$. The perturbation radius is searched from $\{1, 5\}$ and the rate is set to 1. We run each model with 3 random seeds and report the average results of the test set for each method.

4 Results and Discussion

4.1 Overall Results

The overall multi-domain fake news detection results on CrossEval and CovidEval are listed in Table 1 and Table 2, respectively.

Methods	# Para.	Weibo21 [◊]		Twi	tter♦	WeiboC	Covid19 [♦]	Avg.			
		Acc	F1	Acc	F1	Acc	F1	Acc	F1		
Single-domain Fake News Detection											
XLM-RoBERTa	$265.2B \times 3$	87.42	87.39	86.39	86.38	85.42	84.61	86.41	86.13		
TextCNN	$266.1B \times 3$	89.29	89.21	86.47	86.41	83.33	82.35	86.36	85.99		
Multi-domain Fake News	Detection										
XLM-RoBERTa-M	265.2B	88.46	88.33	87.40	87.42	86.11	85.71	87.33	87.15		
EANN	265.2B	87.05	86.94	86.29	86.27	86.11	85.78	86.48	86.33		
MDFEND	272.8B	90.80	90.75	85.77	85.78	88.89	88.15	88.49	88.22		
SAT-TGN	265.3B	90.48	90.48	85.52	85.51	86.81	85.98	87.60	87.32		
StruACL-TGN (ours)	265.3B	91.14	91.09	87.35	87.34	89.81	89.31	89.44	89.25		
StruACL-TGN* (ours)	< 265.3B × 3	91.78	91.78	88.45	88.45	90.97	90.33	90.40	90.19		

Table 1: Experimental results of fake news detection on *CrossEval* benchmark, which involves fake news detection datasets across different social platforms. \diamond refers to content-only datasets without propagation thread. \blacklozenge indicates datasets with both textual content and propagation data. StruACL-TGN and StruACL-TGN* indicate the TGN model trained on the full benchmark and pair-wise benchmark under the proposed StruACL objective, respectively. Detailed results of pair-wise benchmark are listed in Table 3.

Methods	Covid19 [◊]		WeiboCovid19 [♦]		TwitterCovid19♦		Arabic♦		Cantonese♦		Avg.	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
Single-domain Fake News Detection												
XLM-RoBERTa	96.73	96.72	64.58	48.19	62.87	38.60	64.13	39.07	70.30	68.59	71.72	58.23
TextCNN	97.10	97.09	83.33	82.35	67.33	52.06	80.43	75.28	76.66	74.45	80.97	76.25
Multi-domain Fake News Detection												
XLM-RoBERTa-M	96.92	96.90	85.42	84.80	63.37	60.07	77.17	71.52	66.99	61.23	77.97	74.91
EANN	96.40	96.39	86.80	86.40	73.27	66.65	76.09	73.65	63.95	63.73	79.30	77.36
MDFEND	96.40	96.39	84.72	84.03	67.82	59.17	75.54	70.67	72.51	69.58	79.40	75.97
SAT-TGN	96.96	96.95	86.11	85.29	68.32	59.56	82.07	79.85	69.06	67.51	80.50	77.83
StruACL-TGN (ours)	95.69	95.65	86.81	86.21	75.41	72.21	85.33	83.27	71.50	69.55	82.95	81.38

Table 2: Experimental results of fake news detection on *CovidEval* benchmark, which involves fake news detection datasets related to the breaking event COVID-19. \diamond refers to content-only datasets without propagation thread. \blacklozenge indicates datasets with both textual content and propagation data.

						-		-	~			1			
Mathada	Weib	o21 [◊]	WeiboCo	ovid19	A	vg.		Mathada	We	ibo21°	Twitter▼		Avg.		
Methods	Acc	F1	Acc	F1	Acc	F1		wiethous	Acc	F1	Acc	F1	Acc	F1	
XLM-RoBERTa-M	86.17	86.12	87.50	86.85	86.83	86.48	-	XLM-RoBERTa-M	90.74	90.71	86.27	86.26	88.50	88.49	
EANN	80.56	80.09	90.43	90.43	85.49	85.24		EANN	89.08	89.05	86.29	86.29	87.69	87.67	
MDFEND	88.92	88.84	89.58	89.00	89.25	88.92		MDFEND	88.66	88.62	83.69	83.69	86.18	86.16	
SAT-TGN	87.50	87.08	85.28	85.18	86.39	86.13		SAT-TGN	88.82	88.82	82.79	82.70	85.81	85.76	
StruACL-TGN	89.18	89.15	90.97	90.33	90.08	89.74		StruACL-TGN	91.78	91.78	88.45	88.45	90.12	90.12	
w/o StruACL	84.56	84.52	87.73	87.20	86.14	85.86		w/o StruACL	90.74	90.69	87.27	87.27	89.01	88.98	
(a) Detection on Weibo21 and WeiboCovid19.								(b) Detection on Weibo21 and Twitter.							
	WeiboCovid19 [♦] TwitterCovid19		Avg.			Mathada	Weib	Weibo21 [♦]		Covid19 [♦]		g.			
Methods	Acc	F1	Acc	F1	Ad	∞ F	1	Methous	Acc	F1	Acc	F1	Acc	F1	
XLM-RoBERTa-M	86.81	86.16	5 76.24	73.94	81.	52 80.	05 -	XLM-RoBERTa-M	90.74	90.73	97.52	97.52	94.13	94.12	
EANN	87.50	86.56	5 74.26	73.51	80.	88 80.	04	EANN	89.24	89.23	97.34	97.33	93.29	93.28	
MDFEND	82.64	80.57	7 71.78	64.60	77.	21 72.	58	MDFEND	88.46	88.46	97.29	97.28	92.87	92.87	
SAT-TGN	88.19	87.54	4 76.24	73.40	82.	22 80.	47	SAT-TGN	90.95	90.95	97.57	97.56	94.26	94.26	
StruACL-TGN	88.89	88.38	3 78.22	76.40	83.	55 82.	39	StruACL-TGN	92.04	92.01	98.18	98.17	95.11	95.09	
w/o StruACL	86.81	85.98	3 75.74	74.09	81.	27 80.	03	w/o StruACL	90.64	90.63	97.48	97.47	94.06	94.05	
w/o StruACL	86.81	85.98	5 /5.74	- /4.09	81.	27 80.	03	w/o StruACL	90.64	90.63	97.48	97.47	94.06	94.05	

(c) Detection on TwitterCovid19 and WeiboCovid19.

(d) Detection on Weibo21 and Covid19.

Table 3: Multi-domain detection results on pair-wise fake news datasets.

Comparison with Single-domain Fake News Detection Compared with methods in the first block of two tables, our StruACL achieves better performance under lighter network architecture, showing the effectiveness and efficiency of our method. Specifically, StruACL outperforms single-domain detection models by +5.1% and 3.1% in F1 scores on CovidEval and CrossEval, respectively.

399

400

401

402

403

404

405

406

Comparison under Multi-domain Detection Among multi-domain detection methods in the second block of two tables, our StruACL-TGN and StruACL-TGN* achieve the best overall detection performance based on average metrics for accuracy and F1 scores on both benchmarks, showing the superiority of StruACL for multidomain fake news detection. Specifically, for 407

408

409

410

411

412

413

Mathada		Weibo21		Twitter▼ \		W	eiboCo	vid19▼	A	vg.			
	methous		Acc	F1	Acc	F1	Α	cc	F1	Acc	F1		
5	StruACl	Ĺ	91.14	91.09	87.35	87.34	89	.81	89.31	89.44	89.25	-	
	w/o Ad	v	89.93	89.91	86.53	86.50	87	.96	87.41	88.13	87.94		
	w/o Str	uCL	91.21	91.20	87.20	87.20	89	.58	89.13	89.33	89.18		
	w/o Str	uACL	88.30	88.25	85.59	85.57	87	.04	86.55	86.97	86.79		
	w/o Ad	VStruCL	91.23	91.19	87.02	87.00	89	.35	88.77	89.29	88.98	_	
	w/o Ad	VCLS	88.96	88.91	84.29	84.28	88	.19	87.66	87.15	86.95		
(a) Ablation results on <i>CrossEval</i> benchmark.													
Mathada	Covi	d19◊	WeiboO	Covid19•	TwitterCovid1		9 †	Ara	ıbic♦	Cantonese [♦]		Avg.	
Methous	Acc	F1	Acc	F1	Acc	F1		Acc	F1	Acc	F1	Acc	F1
StruACL	95.69	95.65	86.81	86.21	75.41	72.21	l	85.33	83.27	71.50	69.55	82.95	81.38
w/o Adv	96.21	96.21	88.89	88.31	76.24	72.32	2	82.61	78.80	65.76	65.00	81.94	80.13
w/o StruCL	95.61	95.58	86.11	85.84	73.27	70.68	3	84.24	81.29	70.44	68.67	81.93	80.41
w/o StruACL	95.79	95.78	86.81	86.46	73.27	68.76	5	83.33	81.06	68.23	67.16	81.49	79.84
w/o AdvStruCL	97.88	97.88	82.18	81.91	68.32	62.60)	80.62	77.22	70.30	67.62	79.86	77.45

(b) Ablation results on CovidEval benchmark.

66.48

72.28

Table 4: Ablation results of the proposed StruACL on two benchmarks.

CovidEval, compared with the baseline, our Stru-415 ACL achieves 1.46%/1.54% accuracy/F1 scores. 416 For CrossEval, StruACL outperforms the base-417 line about 2.46%/2.47% accuracy and F1 scores. 418 EANN and SAT use adversarial training to learn 419 domain-invariant or structure-invariant features 420 across multi-domain data. MDFEND models 421 domain-specific and domain-shared semantics with 422 complex neural networks. All methods ignore po-423 tential connections between semantics and struc-424 ture. StruACL effectively performs knowledge 425 426 transfer not only across multiple domains but also between semantics and structures, achieving the 427 superior multi-domain detection performance. 428

w/o AdvCLS

96.73

96.72

88.89

88.23

Pair-wise Domain Fake News Detection Ta-429 430 ble 3 shows multi-domain detection results on pairwise datasets where (a) and (b) indicate results 431 trained on two heterogeneous datasets, one with 432 433 propagation and another without; (c) and (d) indicates results trained on two propagation-based 434 datasets and two content-based datasets with a ho-435 mogeneous setting. From results, our proposed 436 StruACL-TGN obtains superior average perfor-437 mance consistently, showing the effectiveness of 438 StruACL on both heterogeneous and homogeneous 439 settings. Additionally, by comparing Table 3 (a) 440 and (b), which display results under two settings 441 trained with Weibo21, we observe that our method 442 443 shows a more significant improvement on Weibo-Covid19, which has less training data compared 444 to Twitter. This suggests our method is capable 445 of better generalization and adaptation to data-446 constrained domains. 447

4.2 Ablation Study

83.70

81.12

67.68

66.27

81.85

79.76

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

476

477

478

479

We further ablate the key components to evaluate the effectiveness of StruACL objective. w/o Adv refers to removing all adversarial perturbations during training. We also remove the perturbation based on structure-aware contrastive learning and crossentropy classification objectives, respectively, denoted as w/o Adv_{StruCL}, and w/o Adv_{CLS}. w/o StruCL indicates removing the structure-aware contrastive learning, ignoring the transfer learning between structure and semantic. w/o StruACL is removing the full StruACL objective. As shown in Table 4, the full model gains the best performance on both benchmarks, compared with the ablated models w/o Adv, w/o StruCL, w/o StruACL. The results demonstrate the effectiveness of each key component for detection. Additionally, for generating adversarial samples, eliminating the guidance of either structure-aware contrastive learning or task prediction gains (i.e., w/o Adv_{StruCL}, and w/o Adv_{CLS}) decreases performance to some extent, demonstrating the effectiveness of both objectives.

4.3 Generalization Evaluation with Data-constrained Conditions

We evaluate the generalization under multi-domain data-constrained conditions. We vary the training set ratios to evaluate detection performance under limited data conditions. Specifically, for a predefined ratio (e.g., 20%), we randomly sampled subsets from the original training sets of all domains. All methods are tested using on the same sampled training subsets to ensure a fair comparison. Fig-



Figure 3: Results against removing domain-specific propagation in the training sets on CrossEval and CovidEval.



Figure 4: Results against different training set sizes. We report the average F1s of datasets on each benchmark.

ure 4 shows results of representative methods and our StruACL on CrossEval and CovidEval across various training set sizes. Our proposed StruACL consistently achieves superior performance across all data-constrained settings, regardless of the training set ratio. This demonstrates the strong generalization capabilities of StruACL in scenarios with limited data. The performance improvement of StruACL is not only attributed to its ability to effectively learn semantic and structural features from multi-domain data but also to its capacity for transferring these learned semantic and structural representations across tasks. These advantages enable StruACL to efficiently utilize limited data and achieve generalized performance in dataconstrained scenarios.

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

4.4 Effect of Training Propagation Structure

We analyze the effect of propagation structures in the training data during transferring between semantic and structure. We remove the propagation structure of the training set on the specific domain. As shown in Figure 3, after removing propagation structures on Twitter and WeiboCovid19, the detection performance on all three datasets declined consistently. This indicates that propagation structures play a critical role in identifying fake news, as they provide complementary signals that enhance semantic analysis. Our StruACL can fully model interactions between semantic and structural features, thereby boosting multi-domain fake news detection. Interestingly, when removing the propagation data of WeiboCovid19, the detection performance on the Twitter declined more significantly compared to the performance drop observed for WeiboCovid19 itself. This may be because that StruACL leverages latent semantic associations related to Weibo, which facilitates the detection on WeiboCovid19 even in the absence of propagation features. In contrast, the performance gains on TwitterCovid19 are primarily driven by the transfer of propagation features from WeiboCovid19. This underscores the critical role of transferable propagation structure features in multi-domain detection.

499

501

503

504

505

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

528

529

530

531

532

533

534

535

536

5 Conclusion

This paper studies a propagation-adaptive multidomain fake news detection paradigm. To achieve this, we develop a novel StruACL-TGN to learn generalized representations from propagationbased and content-only domains. StruACL contrasts semantic and structural representations to preserve and transfer structural knowledge, while introduces propagation-guided adversarial training to enhance the diversity of representations. Experiments on seven datasets show that StruACL-TGN achieves superior multi-domain detection on general and data-constraint settings, proving the effectiveness and generalizability of StruACL.

537 Limitations

The current framework focuses on text-based semantic and propagation data. The study of multimodal inputs, such as images and videos will be left as future work to further enhance the robustness and versatility of fake news detection systems in increasingly complex and dynamic information ecosystems.

References

546

547

548

550

552

553

554

556

557

560

563

565

566

568

569

573

574

575

576

577

579

581

582

584

585

588

- Firoj Alam, Shaden Shaar, Fahim Dalvi, Hassan Sajjad, Alex Nikolov, Hamdy Mubarak, Giovanni Da San Martino, Ahmed Abdelali, Nadir Durrani, Kareem Darwish, et al. 2021. Fighting the covid-19 infodemic: Modeling the perspective of journalists, fact-checkers, social media platforms, policy makers, and the society. In *Findings of the Association* for Computational Linguistics: EMNLP 2021, pages 611–649.
- Tian Bian, Xi Xiao, Tingyang Xu, et al. 2020. Rumor detection on social media with bi-directional graph convolutional networks. In *AAAI*, pages 549–556.
- Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. 2011. Information credibility on twitter. In *WWW*, pages 675–684.
- Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. 2020. Simple and deep graph convolutional networks. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 1725– 1735. PMLR.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In ACL, pages 8440–8451. Association for Computational Linguistics.
- Dou Hu, Lingwei Wei, Wei Zhou, et al. 2021. A rumor detection approach based on multi-relational propagation tree. *Journal of Computer Research and Development*, 58(7):1395–1411.
- S. Mo Jang, Tieming Geng, Jo-Yun Queenie Li, et al. 2018. A computational approach for examining the roots and spreading patterns of fake news: Evolution tree analysis. *Comput. Hum. Behav.*, 84:103–113.
- Debanjana Kar, Mohit Bhardwaj, Suranjana Samanta, and Amar Prakash Azad. 2021. No rumours please! a multi-indic-lingual approach for covid fake-tweet detection. In 2021 grace hopper celebration India (GHCI), pages 1–5. IEEE.
- Hamid Karimi and Jiliang Tang. 2019. Learning hierarchical discourse-level structure for fake news detection. In *NAACL-HLT*, pages 3432–3442.

Liang Ke, Xinyu Chen, Zhipeng Lu, Hanjian Su, and Haizhou Wang. 2020. A novel approach for cantonese rumor detection based on deep neural network. In 2020 IEEE International conference on systems, man, and cybernetics (SMC), pages 1610– 1615. IEEE. 589

590

592

593

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

- Ling Min Serena Khoo, Hai Leong Chieu, Zhong Qian, and Jing Jiang. 2020. Interpretable rumor detection in microblogs by attending to user interactions. In *AAAI*, pages 8783–8790.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. In *EMNLP*, pages 1746–1751. ACL.
- Thomas N. Kipf and Max Welling. 2017. Semisupervised classification with graph convolutional networks. In *ICLR (Poster)*.
- Sumeet Kumar and Kathleen M. Carley. 2019. Tree lstms with convolution units to predict stance and rumor veracity in social media conversations. In *ACL*, pages 5047–5058.
- Jiayang Li, Xuan Feng, Tianlong Gu, and Liang Chang. 2024. Dual-teacher de-biasing distillation framework for multi-domain fake news detection. In *ICDE*, pages 3627–3639. IEEE.
- Chaoqi Liang, Yu Zhang, Xinyuan Li, Jinyu Zhang, and Yongqi Yu. 2022. Fudfend: fuzzy-domain for multidomain fake news detection. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 45–57. Springer.
- Hongzhan Lin, Jing Ma, Liangliang Chen, Zhiwei Yang, Mingfei Cheng, and Chen Guang. 2022. Detect rumors in microblog posts for low-resource domains via adversarial contrastive learning. In *Findings* of the Association for Computational Linguistics: NAACL 2022, pages 2543–2556.
- Hongzhan Lin, Jing Ma, Mingfei Cheng, et al. 2021. Rumor detection on twitter with claim-guided hierarchical graph attention networks. In *EMNLP*, pages 10035–10047.
- Hongzhan Lin, Pengyao Yi, Jing Ma, Haiyun Jiang, Ziyang Luo, Shuming Shi, and Ruifang Liu. 2023. Zero-shot rumor detection with propagation structure via prompt learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 5213–5221.
- Sahil Loomba, Alexandre de Figueiredo, Simon J Piatek, et al. 2021. Measuring the impact of covid-19 vaccine misinformation on vaccination intent in the uk and usa. *Nature human behaviour*, 5:337–348.
- Jing Ma and Wei Gao. 2020. Debunking rumors on twitter with tree transformer. In *COLING*, pages 5455–5466.
- Jing Ma, Wei Gao, Prasenjit Mitra, et al. 2016. Detecting rumors from microblogs with recurrent neural networks. In *IJCAI*, pages 3818–3824.

- 662 664 670 671 672 673 675 679 681 688

- Nikhil Mehta, Maria Leonor Pacheco, and Dan Goldwasser. 2022. Tackling fake news detection by continually improving social context representations using graph neural networks. In ACL, pages 1363-1380.
 - Takeru Miyato, Andrew M. Dai, and Ian J. Goodfellow. 2017. Adversarial training methods for semisupervised text classification. In ICLR (Poster).

Jing Ma, Wei Gao, and Kam-Fai Wong. 2018. Rumor

neural networks. In ACL, pages 1980-1989.

detection on twitter with tree-structured recursive

- Qiong Nan, Juan Cao, Yongchun Zhu, Yanyan Wang, and Jintao Li. 2021. MDFEND: multi-domain fake news detection. In CIKM, pages 3343–3347. ACM.
- Parth Patwa, Shivam Sharma, Srinivas PYKL, Vineeth Guptha, Gitanjali Kumari, Md. Shad Akhtar, Asif Ekbal, Amitava Das, and Tanmoy Chakraborty. 2021. Fighting an infodemic: COVID-19 fake news dataset. In CONSTRAINT@AAAI, volume 1402 of Communications in Computer and Information Science, pages 21-29. Springer.
- Kashyap Popat. 2017. Assessing the credibility of claims on the web. In WWW (Companion Volume), pages 735-739.
- Martin Potthast, Johannes Kiesel, Kevin Reinartz, et al. 2018. A stylometric inquiry into hyperpartisan and fake news. In ACL, pages 231-240.
- Natali Ruchansky, Sungyong Seo, and Yan Liu. 2017. CSI: A hybrid deep model for fake news detection. In CIKM, pages 797–806.
- Michael Sejr Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. In ESWC, volume 10843 of Lecture Notes in Computer Science, pages 593-607.
- Kai Shu, Limeng Cui, Suhang Wang, et al. 2019. defend: Explainable fake news detection. In KDD, pages 395-405.
- Yu Tong, Weihai Lu, Zhe Zhao, Song Lai, and Tong Shi. 2024. MMDFND: multi-modal multi-domain fake news detection. In ACM Multimedia, pages 1178-1186. ACM.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In NIPS, pages 5998-6008.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, et al. 2018. Graph attention networks. In ICLR (Poster).
- Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. The spread of true and false news online. Science, 359(6380):1146-1151.

Yaqing Wang, Fenglong Ma, Zhiwei Jin, Ye Yuan, Guangxu Xun, Kishlay Jha, Lu Su, and Jing Gao. 2018. Eann: Event adversarial neural networks for multi-modal fake news detection. In Proceedings of the 24th acm sigkdd international conference on knowledge discovery & data mining, pages 849–857. 695

696

697

698

699

701

702

704

707

708

709

710

711

712

713

714

715

716

717

718

719

721

722

723

724

726

727

728

729

- Lingwei Wei, Dou Hu, Yantong Lai, et al. 2022a. A unified propagation forest-based framework for fake news detection. In COLING, pages 2769-2779.
- Lingwei Wei, Dou Hu, Wei Zhou, and Songlin Hu. 2024. Transferring structure knowledge: A new task to fake news detection towards cold-start propagation. In ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 8045-8049. IEEE.
- Lingwei Wei, Dou Hu, Wei Zhou, et al. 2021. Towards propagation uncertainty: Edge-enhanced bayesian graph convolutional networks for rumor detection. In ACL, pages 3845-3854.
- Lingwei Wei, Dou Hu, Wei Zhou, et al. 2022b. Uncertainty-aware propagation structure reconstruction for fake news detection. In COLING, pages 2759-2768.
- Ruichao Yang, Xiting Wang, Yiqiao Jin, Chaozhuo Li, Jianxun Lian, and Xing Xie. 2022. Reinforcement subgraph reasoning for fake news detection. In KDD, pages 2253-2262.
- Huaiwen Zhang, Shengsheng Qian, Quan Fang, and Changsheng Xu. 2021. Multimodal disentangled domain adaption for social media event rumor detection. IEEE Trans. Multim., 23:4441-4454.
- Yongchun Zhu, Qiang Sheng, Juan Cao, Qiong Nan, Kai Shu, Minghui Wu, Jindong Wang, and Fuzhen Zhuang. 2023. Memory-guided multi-view multidomain fake news detection. IEEE Trans. Knowl. Data Eng., 35(7):7178-7191.

822

823

824

825

826

827

828

781

Overall of Appendix

731

732

733

734

737

738

740

741

742

743

744

745

747

752

753

755

757

759

764

765

767

772

773

774

776

777

778

780

In the appendix, we will provide related work on fake news detection, and detailed experimental setups.

A Related Work

A.1 Fake News Detection

Fake news detection aims to automatically identify a news piece as real or fake.

Early works on **content-based fake news detection** rely on feature engineering to capture textual characteristics, e.g., topic features (Castillo et al., 2011), writing styles and consistency (Popat, 2017; Potthast et al., 2018). After the emergence of deep learning, some works (Ma et al., 2016; Ruchansky et al., 2017; Karimi and Tang, 2019) applied various neural networks to learn high-level linguistic features from the source news or combing its retweets.

Generally, users on social media share opinions, conjectures and evidence for checking fake news. Through their various interactive behaviors, a propagation tree describing the law of information transmission is formed and plays a significant role in fake news detection. Vosoughi et al. (2018); Jang et al. (2018) have empirically shown that compared with the truth, false news has deeper propagation structures, and reaches a wider audience. To leverage structure properties, propagation-based fake news detection models (Ma et al., 2016; Shu et al., 2019; Khoo et al., 2020) learn the sequential structure features in the propagation trees by RNN-based or attention-based modules. (Shu et al., 2019) jointly learned the sequential effect of comments and correlation between source news and the corresponding comments. To capture structural propagation patterns, (Ma et al., 2016) constructed a tree-structured neural network to model the propagation structure. (Khoo et al., 2020) adopted Transformer (Vaswani et al., 2017) to learn longdistance interactions. Recently, many researchers (Bian et al., 2020; Hu et al., 2021; Lin et al., 2021; Wei et al., 2021, 2022b; Mehta et al., 2022; Yang et al., 2022) regard the propagation tree as a graph, and employ various graph-based models (Kipf and Welling, 2017; Schlichtkrull et al., 2018; Chen et al., 2020; Velickovic et al., 2018) to capture topological structure features for detection. applied two graph convolutional networks (GCNs) (Kipf and Welling, 2017) to learn structural patterns from two distinct directed graphs. (Hu et al., 2021; Lin

et al., 2021) further lored multi-relational interactions in the propagation graph. Wei et al. (2024) study cold-start propagation and lore transferable features from samples with propagation for improving detection of content-only samples.

A.2 Fake News Detection across Multiple Domains

In real-world scenarios, fake news typically originates and propagate across various domains or platforms, due to real-time events, social trends, and other factors. Thus, multi-domain fake news detection has draw significant attention.

Most works aims to study domain-shared (Zhu et al., 2023; Liang et al., 2022) and domaininvariant semantic features (Wang et al., 2018; Zhang et al., 2021; Li et al., 2024) for detecting fake news across multiple domains. For example, Wang et al. (2018) learn event-invariant representations for multi-domain detection via considering the effect of event diversity. Nan et al. (2021) utilize domain gate to alleviate the domain shift issue for aggregation of multi-domain representations. Li et al. (2024) study the unbalanced multi-domain data issue and leverage two teacher models to mitigate the domain bias via knowledge distillation. (Zhu et al., 2023) introduce a domain adapter to extract domain-shared features from similar domains for fake news detection. Liang et al. (2022) design a fuzzy domain label to lore multidomain knowledge. Tong et al. (2024) design a progressive hierarchical extraction network to achieve domain-adaptive domain-related knowledge extraction. Most of the above multi-domain methods focus on the news content across different domains, ignoring potential shared propagation structures for detection.

B Datasets

All experiments are conducted on seven fake news detection datasets. Specifically, **Weibo21** (Nan et al., 2021) collects Chinese tweets without propagation data on Sina Weibo platform ranging from 2010-12-15 to 2021-03-31. Regarding to the breaking event COVID-19 pandemic, **Covid19** (Patwa et al., 2021) collects English textual tweets related to the topic of COVID-19 from from public fact verification websites and social media (e.g., Facebook and Twitter¹). **TwitterCovid19** Kar et al. (2021); Lin et al. (2022) collects English textual tweets

¹Renamed X in 2023.

Dataset	Prop.	Resource	Lang.	Event	# Train	# Valid	# Test	# Total
Weibo21	×	Weibo	CN	Hybrid	5,751	1,918	1,923	9,592
Covid19	×	Hybrid	EN	COVID19	6,420	2,140	2,140	10,700
Twitter	\checkmark	Twitter	EN	Hybrid	3,109	777	3,888	7,774
WeiboCovid19	\checkmark	Weibo	CN	COVID19	163	40	208	411
TwitterCovid19	\checkmark	Twitter	EN	COVID19	159	39	202	400
Arabic	\checkmark	Twitter	AR	COVID19	136	36	184	356
Cantonese	\checkmark	Twitter	YUE	COVID19	577	143	724	1,444

Table 5: Statistics of 7 datasets for fake news detection. Prop. refers to whether the dataset contains propagation data. Lang. indicates language used in the dataset where CN, EN, AR, and YUE represent Chinese, English, Arabic, and Cantonese, respectively. Event summarizes the types of social events collected in the dataset. # Total refers to the total number of samples in each dataset.

and the corresponding propagation from Twitter. 829 WeiboCovid19 (Lin et al., 2022) contains relevant 830 Chinese textual claims and the propagation thread 831 on Sina Weibo and X platform. Arabic and Can-832 833 tonese, originally collected by Alam et al. (2021) and Ke et al. (2020), contain textual claims in Ara-834 bic and Cantonese, respectively. Lin et al. (2023) 835 further collect the propagation thread of each claim 836 tweets on both datasets. The statistics of the above 837 838 datasets are shown in Table 5.