Context-Aware Deepfake Detection for Securing AI-Driven Financial Transactions

Changcun Liu[†], Guisheng Zhang[†], Siyou Guo, Qilei Li, Gwanggil Jeon, Mingliang Gao*

Abstract—The rapid advancement of deepfake technology has threatened the community's sense of security, particularly in 2 the context of face-based payment systems. Thus, deepfake 3 detection has emerged as a critical issue demanding immediate 4 attention. However, the generalization performance of existing 5 detection models is limited as they are overly reliant on specific forged features while ignoring the common forged features. To address this problem, we introduce the Context-Aware Decou-8 pling Network (CADNet) for deepfake detection. Specifically, a Context Self-Calibration (CSC) module is constructed to guide 10 the network to focus on local forged regions. It enlarges possible 11 regions to increase the likelihood of forgery cues. Meanwhile, a 12 Frequency Domain Decoupling (FDD) module is introduced to 13 14 extract and fuse different frequency components. It realizes the collaborative representation optimization of global semantics and 15 local details. The experimental results prove that the proposed 16 model exhibits strong generalization capability across multiple 17 standard datasets. It achieves average AUC values of 98.64% for 18 in-domain evaluation and 75.52% for cross-dataset generaliza-19 tion. 20

Index Terms—Deepfake detection, Facial payment, Generalization, Frequency domain decoupling, Context self-calibration.

I. INTRODUCTION

24

Ith the rapid development of the internet and the 25 flourishing digital economy, digital finance has become 26 a pivotal driver of economic growth. In this context, it is 27 increasingly vital to ensure secure transactions, particularly in 28 facial payment systems. However, with the rapid development 29 of Artificial Intelligence Generated Content (AIGC), the mis-30 use of forged facial images has significantly increased in facial 31 recognition payment systems. As digital finance expands, the 32 risks posed by fraudulent facial images in the financial sector 33 have become increasingly acute. To address this challenge, 34 the academic community has introduced deepfake detection 35 methods and conducted systematic and extensive research on 36 deepfake detection [1]-[3]. 37

Existing deepfake detection methods typically perform well when training and testing data are generated using the same

Gwanggil Jeon is with the Department of Embedded Systems Engineering, Incheon National University, Incheon, 22012, South Korea (Email:ggjeon@gmail.com).

Qilei Li is with the School of Electronic Engineering and Computer Science, Queen Mary University of London, London, E1 4NS, United Kingdom (Email: q.li@qmul.ac.uk).

* Mingliang Gao is the corresponding author.

[†] Changcun Liu and Guisheng Zhang contributed equally to this work.

deepfake techniques [4]. However, in real-world applications, 40 testing data may be generated using unknown methods, and 41 thus, the generalization performance is degraded. To improve 42 the generalization performance of deepfake detection, some 43 researchers have made some attempts in the spatial and fre-44 quency domains. Yan et al. [5] improved the generalization by 45 mining the common features in different forgery techniques. 46 Liu et al. [6] introduced a spatial-phase shallow learning 47 method to improve the generalization ability of forgery face 48 detection. Qian et al. [7] developed a frequency-based face 49 forgery detection network to enhance the generalization capa-50 bility. Luo et al. [8] improves the generalization of forgery 51 detection systems by exploiting high-frequency features. Al-52 though these methods aim to improve the generalization ability 53 of detection networks, their performance remains limited. 54 This is because deepfake detection methods overfit to specific 55 prominent regions of the deepfake images [9]. 56

To improve the generalization of the deepfake detection 57 model, we propose the CADNet for deepfake detection. The 58 CADNet mainly comprises two key modules, namely Context 59 Self-Calibration (CSC) and Frequency Domain Decoupling 60 (FDD). The CSC module retains the original input character-61 istics and extracts local features. It also prompts the network 62 to explore more forged content by examining as many regions 63 as possible. Additionally, the FDD module extracts and fuses 64 the low-frequency and three high-frequency components. The 65 low-frequency component is used to capture the global pattern, 66 and the three directional high-frequency components are used 67 to retain detailed spatial information. In the process of high 68 and low frequency feature processing, the EffConv block is 69 designed to eliminate redundant feature representations. In 70 sum, the contributions of this work are three-fold: 71

- A CADNet model is proposed to enhance the generalization and accuracy of the deepfake detection.
- A CSC module is introduced to guide the network to focus on the forged regions and assess specific non-critical regions effectively.
- An FDD module is developed to extract and fuse lowfrequency and high-frequency components to enhance feature extraction.

The rest of the paper is structured as follows: Section II presents the related work. Section III illustrates the proposed method in detail. Section IV analyses the experimental results. The paper is concluded in Section V.

74 75 76

77

78

79

72

Changcun Liu, Guisheng Zhang, Siyou Guo, and Mingliang Gao are with the School of Electrical and Electronic Engineering, Shandong University of Technology, Zibo, 255000, China. (Email: 24404020554@stumail.sdut.edu.cn, 22504030001@stumail.sdut.edu.cn, 23504030565@stumail.sdut.edu.cn, and mlgao@sdut.edu.cn.)

85 A. Deepfake detection

Deepfake is a generative technique that uses methods such 86 as Generative Adversarial Network (GAN) and Diffusion 87 model to realistically synthesize or manipulate facial images. 88 In digital finance and other sensitive domains, the misuse of 89 deepfake technology can lead to significant property dam-90 age [10]. Therefore, detecting fake facial images is essential 91 to prevent fraudulent activities and mitigate potential risks. 92 Various approaches have been proposed for deepfake detec-93 tion [11]–[16]. In the early stages, deepfake detection methods 94 focus mainly on identifying abnormal features [17] or forged 95 traces [18] in facial images to determine their authenticity. 96 However, with the emergence of various new deepfake genera-97 tion models, these methods exhibit poor detection performance 98 on advanced deepfake techniques, and these results in poor 99 generalization performance. 100

To enhance the generalization capability of deepfake de-101 tection methods, many deep learning techniques have been 102 extensively adopted [19], benefiting from their strong feature 103 modeling capabilities. Cao et al. [20] proposed an end-to-104 end reconstruction-classification learning framework, termed 105 RECCE, to improve generalization in deepfake detection. 106 It integrates multi-scale graph reasoning and reconstruction-107 guided attention mechanisms to detect forged images with 108 unknown patterns. Ni et al. [21] proposed a consistent rep-109 resentation learning model to improve generalization in de-110 tecting unknown forgery patterns. It introduces a consistent 111 representation learning framework that explicitly constrains 112 feature consistency across different data augmentations using a 113 cosine similarity loss. Dang et al. [22] introduced an attention-114 driven deepfake detection model to identify manipulated facial 115 images and locate tampered facial regions. However, these 116 methods tend to overly rely on specific regional features and 117 ignore the common forged features. Thus, their generalization 118 capabilities are limited. To this aim, we propose the CAD-119 Net through context-dependent insights and decoupled feature 120 analysis. 121

122 B. Global context attention

The global context attention is widely used in deepfake de-123 tection for capturing global information efficiently. It bridges 124 the gap between local and global features and enhances model 125 adaptability and performance. Li et al. [23] proposed a long-126 distance attention mechanism that integrates global informa-127 tion to enhance the comprehensiveness of feature representa-128 tion. Wang et al. [24] introduced the multi-domain attention 129 mapping learning mechanism to enhance the adaptability of 130 the model to diverse deepfakes. Additionally, Das et al. [25] 131 developed a gated context attention mechanism to filter and 132 aggregate relevant contextual information from coarse feature 133 maps. In this work, we build a Context Self-Calibration (CSC) 134 module to guide the network to focus on the forged regions 135 and assess specific non-critical regions. 136

C. Frequency domain analysis

Spatial domain processing methods directly analyze the raw 138 pixels to detect forgery. They overlook frequency domain 139 anomalies generated during the forgery process. In contrast, 140 the frequency domain processing methods detect forgery 141 through spectral transformations. They focus on local fre-142 quency domain features of the image, such as high-frequency 143 components like edges and textures, and low-frequency com-144 ponents like smooth regions. Thus, the frequency domain 145 approaches are more suitable for capturing pixel-level artifacts 146 or frequency anomalies generated during the forgery process. 147

Wavelet transform is a common tool to capture features by 148 multi-resolution analysis in forgery detection. Wang et al. [26] 149 utilized wavelet transform to derive frequency-domain features 150 and improve the generalization of the deepfake detection 151 model across datasets. Gao et al. [27] employed the discrete 152 wavelet transform to enhance global high-frequency features 153 in complex forgery detection. Liu et al. [28] proposed a multi-154 scale wavelet transformer architecture for deepfake detection. 155 Wolter et al. [29] proposed a wavelet packet-based approach 156 for deepfake detection. In this work, we propose a Fre-157 quency Domain Decoupling module to decouple and integrate 158 low-frequency and high-frequency components. The proposed 159 CADNet enhances the model's generalization capability. It 160 uses the integration of multi-level high-frequency features, 161 spatial features, and a frequency-guided attention mechanism. 162

III. PROPOSED MODEL

A. Overview

The schematic of the CADNet is shown in Fig. 1. The 165 image pair is processed by two independent encoders. A face 166 encoder extracts the initial face features, while a background 167 encoder extracts the initial background features. The CSC 168 module enhances the face and background features and derives 169 F_f and F_b (or R_f and R_b) by emphasizing local features 170 and more feature regions. Next, the initial specific features 171 and common features are obtained by separating the facial 172 features through convolutional layers. Subsequently, the initial 173 specific features and common features are enhanced using the 174 FDD module to derive F_s and F_c (R_s and R_c). The common 175 head is designed to assess the authenticity of the images 176 and enhance generalization to unseen forgery techniques. The 177 specific head is a multi-class classification that distinguishes 178 between various forgery techniques. 179

In addition, the reconstruction module is divided into the 180 self-reconstruction part and the cross-reconstruction part. The 181 self-reconstruction part is to ensure that the decoupled fea-182 tures are consistent with the original image, and the cross-183 reconstruction part is to verify the independence of the decou-184 pled features. The background and face features are processed 185 using three Adaptive Instance Normalization (AdaIN) lay-186 ers [30], as well as four convolutional and four upsampling op-187 erations. Then, the reassembled features are processed by the 188 decoder to reconstruct facial images. The reconstructed images 189 are used for reconstruction loss calculation. The reconstruction 190 loss function facilitates enhanced feature decoupling between 191 background and facial features through image reconstruction. 192

137

163

223



Fig. 1. The overview framework of the CADNet model for deepfake detection.

193 B. Context self-calibration module

To enhance the model's ability to capture local key regions and global contextual dependencies in facial images, a Context Self-Calibration (CSC) module is built. The schematic of the CSC module is shown in Fig. 2.



Fig. 2. The overview of the CSC module. H refers to horizontal pooling. V refers to vertical pooling.

Specifically, the residual feature information is compressed 198 along the vertical axis using horizontal pooling. It extracts 199 global information from each column and generates horizon-200 201 tal global context features. Similarly, feature information is compressed along the horizontal axis using vertical pooling. 202 Subsequently, the horizontal and vertical pooling results are 203 combined to generate a fused global context feature map. The 204 fusion formula is represented as: 205

$$R_attention = HP(F) \oplus VP(F), \tag{1}$$

where F is the input feature. HP(F) and VP(F) denote the horizontal and vertical axial vector, respectively. \oplus represents the broadcast addition operation.

The map $R_{attention}$ is processed through a convolutional layer. Meanwhile, a 3×3 convolution is performed on the input feature F. Then, the results of these two operations are multiplied element-wise to generate F'. The feature F' can be described as:

$$F' = Conv(R_attention) \otimes Conv_{3\times3}(F), \qquad (2)$$

where $Conv_{3\times3}(\cdot)$ is a 3×3 convolution. \otimes represents the element-wise multiplication. After that, the feature F' is passed through Batch Normalization (BN) and a Multilayer Perceptron (MLP) to obtain the global features. Finally, the model adds the global features and the input features to produce the final output. The process is formulated as: 219

$$F_{csc} = \mathrm{MLP}(\mathrm{BN}(F')) + F, \tag{3}$$

where F_{csc} is the output feature of the CSC module. MLP(·) 220 and BN(·) represent the multilayer perceptron and batch 221 normalization, respectively. 222

C. Frequency domain decoupling module

The schematic of the FDD module is shown in Fig. 3. 224 Within the FDD module, the feature map undergoes decomposition through the Haar wavelet transform. This process produces a low-frequency component along with high-frequency components in horizontal, vertical, and diagonal orientations. 228 The Haar wavelet transform formula is expressed as: 229

$$l_r(i,j) = \frac{F_f(i,2j-1) + F_f(i,2j)}{2},$$

$$d_r(i,j) = F_f(2i-1,j) - F_f(2i,j),$$
(4)

where *i* and *j* are the 2D coordinate indices, and they are employed to locate the spatial position of the feature explicitly. $l_r(i, j)$ is a low-frequency component. F_f is the input feature. $d_r(i, j)$ is a diagonal high-frequency component. 233

Then, the low-frequency feature and three high-frequency ²³⁴ features can be obtained as: ²³⁵

$$l(i,j) = \frac{l_r(2i-1,j) + l_r(2i,j)}{2},$$

$$h_L(i,j) = l_r(2i-1,j) - l_r(2i,j),$$

$$h_V(i,j) = d_r(2i-1,j) + d_r(2i,j),$$

$$h_D(i,j) = d_r(2i-1,j) - d_r(2i,j),$$

(5)

where l(i, j) denotes the initial low frequency component. ²³⁶ $h_L(i, j)$ and $h_V(i, j)$ represent the horizontal high-frequency ²³⁷ and vertical high-frequency components, respectively. $h_D(i, j)$ ²³⁸



Fig. 3. The overview of the FDD module. l and h refer to the low-frequency and high-frequency features. h_L and h_V denote the horizontal high-frequency and vertical high-frequency features, respectively. h_D refers to the diagonal high-frequency feature.

(6)

is the diagonal high-frequency component. The three high-239 frequency features are concatenated to generate the ini-240 tial high-frequency feature. Subsequently, the initial high-241 frequency features are refined through convolution normaliza-242 tion and an EffConv block to produce the final high-frequency 243 features. The initial low-frequency features are processed 244 through a convolutional operation followed by an EffConv 245 block to obtain the final low-frequency representations. 246

The final low-frequency and high-frequency components are 247 formulated as: 248

$$L = \text{EffConv}(\text{BN}(Conv_{1\times 1}(l))),$$

$$H = \text{EffConv}(\text{BN}(Conv_{1\times 1}(Concat(h_L, h_V, h_D)))),$$

where L and H are the final low-frequency component and 249 high-frequency component, respectively. EffConv(\cdot) denotes 250 the efficient convolution. $BN(\cdot)$ represents the batch nor-251 malization. $Conv_{1\times 1}(\cdot)$ is the convolution kernel with a size 252 of 1×1 . Concat(·) represents the concatenation operation 253 applied to three high-frequency features. 254

By employing the FDD module, the proposed method main-255 tains high-frequency detail representation while simultane-256 ously integrating low-frequency global structural information. 257

D. Loss function 258

The loss function is a combination of classification loss, 259 contrastive regularization loss, and reconstruction loss. 260

1) Classification loss: The classification loss consists of 261 two components: the binary classification loss and the Specific 262 forgery classification loss. Binary classification loss (\mathcal{L}_{ce}^{c}): It 263 supervises the model to learn common forgery features across 264 different methods. The loss \mathcal{L}_{ce}^{c} is denoted as: 265

$$\mathcal{L}_{ce}^{c} = \mathcal{L}_{ce}\left(\mathcal{H}_{c}\left(A_{c}\right), y_{i}\right),\tag{7}$$

where $\mathcal{H}_{c}(\cdot)$ is the head for common forgery features. A_{c} is the 266 common fingerprint, and $y_i \in \{\text{fake, real}\}$. Specific forgery 267 **classification loss** (\mathcal{L}_{ce}^{s}): It is employed to identify the specific 268 forgery method. The loss is formulated as: 269

$$\mathcal{L}_{ce}^{s} = \mathcal{L}_{ce}\left(\mathcal{H}_{s}\left(A_{s}\right), y_{i}^{\prime}\right),\tag{8}$$

where $\mathcal{H}_s(\cdot)$ is the head for specific forgery features. A_s is 270 the specific fingerprint, and $y'_i \in \{\text{fake}, \text{GAN}_1, \text{GAN}_2, \dots\}$. 271

2) Contrastive regularization loss: The contrastive regular-272 ization loss \mathcal{L}_{con} is formulated as: 273

$$\mathcal{L}_{\text{con}} = \max\left(\|x_A - x_P\|_2 - \|x_A - x_N\|_2 + \alpha, 0\right), \qquad (9)$$

where α is a margin hyperparameter. x_A is the anchor. x_N is 274 the dissimilar counterpart. x_P is the dissimilar counterpart. It 275 ensures a generalizable representation by distinguishing real 276 and fake images and method-specific forgeries. 277

3) Reconstruction loss: The reconstruction loss includes 278 self-reconstruction loss and cross-reconstruction loss. The self-279 reconstruction loss ensures that the reconstructed image aligns 280 with the original by combining the background and face from 281 the same source. The cross-reconstruction loss verifies feature 282 independence by combining a real background with a forged 283 face and vice versa. The final reconstruction loss is expressed 284 as: 285

$$\mathcal{L}_{\text{rec}}^{s} = \|F - D(F_{b} + F_{s+c})\|_{1} + \|R - D(R_{b} + R_{s+c})\|_{1},$$

$$\mathcal{L}_{\text{rec}}^{c} = \|F - D(F_{b} + R_{s+c})\|_{1} + \|R - D(R_{b} + R_{s+c})\|_{1},$$

$$\mathcal{L}_{\text{rec}} = \mathcal{L}_{\text{rec}}^{s} + \mathcal{L}_{\text{rec}}^{c},$$
(10)

where F and R represent the fake and real images. $D(F_b +$ 286 F_{s+c}) and $D(R_b + R_{s+c})$ represent the self-reconstructed 287 images processed by the background decoder and the face 288 decoder. $D(F_b+R_{s+c})$ and $D(R_b+F_{s+c})$ represent the cross-289 reconstructed images processed by the background decoder 290 and the face decoder. $\mathcal{L}^{s}_{rec}, \ \mathcal{L}^{c}_{rec},$ and \mathcal{L}_{rec} correspond to 29 self-reconstruction loss, cross-reconstruction loss, and final 292 reconstruction loss. 293

4) Final loss: The final loss is to balance classification, 294 contrastive, and reconstruction objectives to improve feature 295 disentanglement and detection accuracy. It is obtained as: 296

$$\mathcal{L} = \mathcal{L}_{ce}^c + \lambda_1 \mathcal{L}_{ce}^s + \lambda_2 \mathcal{L}_{con} + \lambda_3 \mathcal{L}_{rec}, \qquad (11)$$

where λ_1 , λ_2 , and λ_3 are weights for specific loss components 297 and their values strictly follow the configuration established in 298 the baseline Uncovering Common Features (UCF) [5].

IV. EXPERIMENTAL RESULTS AND ANALYSIS

A. Datasets

Experiments are conducted on five deepfake datasets, 302 namely FaceForensics++ (FF++) [31], Celeb-DF-v1 [32], 303 Celeb-DF-v2 [32], DeepFake Detection Challenge Pre-304 view (DFDCP) [33] and FaceShifter [34]. The pre-processing 305 for these datasets follows the criteria in work [35]. 306

299

300

 TABLE I

 Comparative results on intra-dataset validation. The training dataset is the FF++ dataset, and the metric is AUC. The best results are highlighted in **Bold**.

Methods	FF++	DF	F2F	FS	NT	Avg.
Meso4 [36]	0.6077	0.6771	0.6170	0.5946	0.5701	0.6133
MesoIncep [36]	0.7583	0.8542	0.8087	0.7421	0.6517	0.7630
CNN-Aug [37]	0.8493	0.9048	0.8788	0.9026	0.7313	0.8534
Xception [31]	0.9637	0.9799	0.9785	0.9833	0.9385	0.9688
EfficientB4 [38]	0.9567	0.9757	0.9758	0.9797	0.9308	0.9637
Capsule [39]	0.8421	0.8669	0.8634	0.8734	0.7804	0.8452
FWA [40]	0.8765	0.9210	0.9000	0.8843	0.8120	0.8788
X-ray [41]	0.9592	0.9794	0.9872	0.9871	0.9290	0.9684
FFD [22]	0.9624	0.9803	0.9784	0.9853	0.9306	0.9674
CORE [21]	0.9638	0.9787	0.9803	0.9823	0.9339	0.9678
Recce [20]	0.9621	0.9797	0.9779	0.9785	0.9357	0.9668
UCF [5]	0.9705	0.9883	0.9840	0.9896	0.9441	0.9753
F3Net [7]	0.9635	0.9793	0.9796	0.9844	0.9354	0.9684
SPSL [6]	0.9610	0.9781	0.9754	0.9829	0.9299	0.9655
SRM [8]	0.9576	0.9733	0.9696	0.9740	0.9295	0.9609
CADNet (Ours)	0.9854	0.9940	0.9929	0.9933	0.9663	0.9864

The FF++ dataset is a comprehensive dataset for training 307 and evaluating deep learning models in facial manipulation 308 detection. With over 500,000 forged images, it surpasses sim-309 ilar datasets and aids in developing robust detection models. 310 The dataset includes four major forgery techniques: Deep-311 Fakes (DF), Face2Face (F2F), FaceSwap (FS), and Neural-312 Textures (NT), which reflect the diverse and advanced trends 313 in facial forgery technology. 314

The Celeb-DF-v1 dataset is a foundational deepfake dataset designed to advance research in detecting manipulated media content. It contains 795 real videos featuring 59 celebrity subjects and 5,639 high-quality deepfake videos. The dataset incorporates early deepfake generation techniques, such as neural network-driven facial synthesis and blending.

The Celeb-DF-v2 dataset does not contain artifacts like unnatural lip movements and geometric distortions. It provides high-resolution videos with an equal mix of real and fake content generated using advanced techniques, which increases the challenge for detection models. As a benchmark dataset, it plays a key role in advancing deepfake detection methods.

The DFDCP dataset provides high-quality real and forged videos that use diverse facial manipulation techniques. It emphasizes realistic scenarios with variations in lighting, angles, and resolutions. This makes it ideal for training and evaluating detection models and supports the development of robust deepfake detection systems.

The FaceShifter dataset is an advanced face-swapping dataset containing highly realistic facial forgeries with minimal artifacts. Its two-stage architecture preserves facial identity and ensures consistency in lighting and background. Using a feature-based generator and refinement network, it delivers seamless results and sets a new benchmark in deepfake detection research.

340 B. Evaluation metrics

The proposed network is evaluated using four standard binary classification metrics, namely Area Under the Curve (AUC), Accuracy (ACC), Average Precision (AP), and Equal Error Rate (EER). The AUC quantifies the ability of the

Input F_b F_f F_c F_s Image: Second seco

Fig. 4. Visualizations of extracted intermediate feature maps.

model to differentiate classes and indicates the probability of 345 correct classification. The ACC represents the model's overall 346 performance in the sample and measures the proportion of 347 samples correctly classified by the model. The AP evaluates 348 the balance between the accuracy and the recall rate of 349 the model at different thresholds. The EER represents the 350 threshold at which the false acceptance rate and false rejection 351 rate converge. 352

C. Implementation details

In this work, we used the Adam [42] optimization algorithm 354 with a learning rate of 2×10^{-4} . The batch size was set to 32. 355 For the Xception [31] detector, we took their official model 356 and initialized the parameters with pre-training on ImageNet. 357 The margin α in Eq. (9) is set to 3. In the final loss function 358 in Eq. (11), the λ_1 , λ_2 , and λ_3 were set as 0.1, 0.05, and 0.3. 359 All experiments were implemented in PyTorch [43] with an 360 NVIDIA GeForce RTX 3090 Ti GPU. 36

D. Experimental results

To maintain strict experimental controls, neither data augmentation techniques nor external training data were incorporated. The model was exclusively trained on the FF++ dataset without additional synthetic or manipulated samples. 366

1) Intra-dataset validation: To conduct intra-testing, all 367 methods are trained on the FF++ dataset and evaluated on 368 both the FF++ dataset and its subsets. The comparative results 369 between different methods are presented in Table I. Compared 370 to other methods, the proposed method achieves superior 371 performance on the FF++ dataset and its subsets. Additionally, 372 when compared to the baseline method, UCF, the CADNet 373 improves the average AUC by 1.11%. Meanwhile, the feature 374 extraction processes of the proposed CADNet on the FF++ 375 dataset are shown in Fig. 4. It shows that the proposed 376 CADNet has the ability to capture the common features of 377 fake images. 378

353



Fig. 5. Visualization of qualitative results on five datasets. The ground truth (GT) and the prediction (Pred) are highlighted in red and blue, respectively.



Fig. 6. Attention visualizations of the proposed CADNet and the baseline (UCF). Compared with the baseline, the proposed CADNet focuses on more facial regions.

2) Cross-dataset validation: To evaluate the cross-domain 379 generalization capability of the proposed model, the cross-380 dataset validation was performed in four unseen out-of-381 distribution datasets. The framework was trained exclusively 382 on the FF++ dataset and subsequently evaluated across 383 four distinct deepfake datasets: Celeb-DF-v1, Celeb-DF-v2, 384 DFDCP, and FaceShifter. The evaluation results are presented 385 in Table II. Compared to other methods, the proposed model 386 performs best on Celeb-DF-v2, DFDCP, and FaceShifter. It 387 388 achieves the second-best result on Celeb-DF-v1, only behind SPSL [6]. 389

390 3) Validation with baseline: To demonstrate the superiority
of the proposed method, we compared the proposed CADNet
with the baseline method, UCF [5]. Comparative results are
presented in Table III. The CADNet outperforms the baseline
models across all metrics on the FF++, DF, F2F, and FS
datasets. On the NT dataset, CADNet surpasses the baseline
in terms of AP and EER.

The subjective results of the CADNet are shown in Fig. 5. It 397 demonstrates the ability of the proposed method to determine 398 the authenticity of images across multiple datasets. Mean-399 while, the visualization of the attention map is presented in 400 Fig. 6. Compared to the UCF, the proposed network can focus 401 on a larger range of facial regions in various datasets. Thus, 402 the proposed CADNet is more capable of detecting forgery 403 traces than the UCF model. 404

 TABLE II

 COMPARATIVE RESULTS ON CROSS-DATASET VALIDATION. THESE

 MODELS ARE TRAINED ON THE FF++ DATASET AND TESTED ON THE

 CELEB-DF-v1, CELEB-DF-v2, DFDCP, AND FACESHIFTER,

 RESPECTIVELY. THE EVALUATION METRIC IS AUC AND AVERAGE AUC

 (AVG.). THE BEST RESULTS ARE HIGHLIGHTED IN BOLD.

Methods	Celeb-DF-v1	Celeb-DF-v2	DFDCP	FaceShifter	Avg.
Meso4 [36]	0.7358	0.6091	0.5994	0.5660	0.6276
MesoIncep [36]	0.7366	0.6966	0.7561	0.6438	0.7083
CNN-Aug [37]	0.7420	0.7027	0.6170	0.5985	0.6651
Xception [31]	0.7794	0.7365	0.7374	0.6249	0.7196
EfficientB4 [38]	0.7909	0.7487	0.7283	0.6162	0.7210
Capsule [39]	0.7909	0.7472	0.6568	0.6465	0.7104
FWA [40]	0.7897	0.6680	0.6375	0.5551	0.6626
X-ray [41]	0.7093	0.6786	0.6942	0.6553	0.6844
FFD [22]	0.7840	0.7435	0.7426	0.6056	0.7189
CORE [21]	0.7798	0.7428	0.7341	0.6032	0.7150
Recce [20]	0.7677	0.7319	0.7419	0.6095	0.7128
UCF [5]	0.7696	0.7391	0.7594	0.6462	0.7286
F3Net [7]	0.7769	0.7352	0.7354	0.5914	0.7097
SPSL [6]	0.8150	0.7650	0.7408	0.6437	0.7411
SRM [8]	0.7926	0.7552	0.7408	0.6014	0.7225
CADNet (Ours)	0.7986	0.7655	0.7683	0.6882	0.7552

E. Analysis of robustness

Considering the ubiquity of image processing, we investi-406 gate the performance under several perturbations, including 407 image compression, Gaussian blur, CLAHE contrast, satu-408 ration, and pixelation. The training and testing dataset is 409 the FF++ dataset. In the details of robustness analysis, we 410 introduce the five standardized perturbation methods. Image 411 compression: The input images were subjected to lossy 412 compression at a quality factor of 70%. Gaussian blur: A 413

TABLE III

COMPARATIVE RESULTS OF THE PROPOSED CADNET AND THE BASELINE (UCF). THE TRAINING DATASET IS THE FF++ DATASET, AND THE METRICS ARE ACC, AP, AUC, AND EER. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD.

	FF++ DF				F2F			FS			NT									
	ACC↑	AP↑	AUC↑	EER↓	ACC↑	AP↑	AUC↑	EER↓	ACC↑	AP↑	AUC↑	EER↓	ACC↑	AP↑	AUC↑	EER↓	ACC↑	AP↑	AUC↑	EER↓
Baseline Ours	0.9445 0.9479	0.9960 0.9965	0.9705 0.9854	0.0589 0.0576	0.9555 0.9588	0.9943 0.9951	0.9883 0.9940	0.0355 0.0330	0.9445 0.9600	0.9924 0.9948	0.9840 0.9929	0.0388 0.0299	0.9049 0.9702	0.9698 0.9948	0.9896 0.9933	0.0922 0.0909	0.9049 0.8798	0.9698 0.9730	0.9441 0.9663	0.0922 0.0909

TABLE IV ROBUSTNESS EVALUATION ON FF++ DATASET. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD.

Methods	Compress	Blur	Contrast	Saturate	Pixelate
Xception [31]	0.5956	0.6176	0.6110	0.5992	0.5996
CORE [21]	0.5508	0.5136	0.5689	0.5365	0.5372
Recce [20]	0.5280	0.5033	0.5472	0.5354	0.5361
UCF [5]	0.5925	0.6247	0.6187	0.5931	0.6055
SPSL [6]	0.5969	0.6280	0.6334	0.6288	0.6284
SRM [8]	0.5583	0.5487	0.5840	0.5672	0.5691
CADNet (Ours)	0.6335	0.6436	0.6673	0.6533	0.6552

 15×15 Gaussian kernel with standard deviation (σ =1.1) was 414 applied for image blurring. CLAHE contrast: Adaptive 415 histogram equalization was performed using an 8×8 grid 416 pattern with a clip limit threshold set to 2.0 for local contrast 417 restriction. Saturation: In the HSV color space, the satu-418 ration channel was amplified by a factor of 1.1. Pixelation: 419 Images were first downsampled to 256×256 resolution and 420 subsequently upsampled to their original dimensions using 421 interpolation. Comparative results are shown in Table IV. 422 The proposed CADNet achieves the best results in image 423 compression, Gaussian blur, CLAHE contrast, saturation, and 424 pixelation operations. It shows that CADNet outperforms other 425 SOTA models when it is subjected to the aforementioned 426 disturbances. 427

F. Ablation study 428

To assess the effectiveness of the core components, ablation 429 experiments were conducted on the FF++ dataset. The results 430 are shown in Table V. In ablation studies, the Baseline is the 431 UCF. It proves that when the CSC and FDD modules are 432 equipped individually, the performance in AUC, ACC, AP, 433 and EER can be improved steadily. When the two modules 434 are equipped simultaneously, AUC, ACC, and AP metrics are 435 further enhanced. The value of EER is slightly higher than the 436 "Baseline+CSC" but it is better than the Baseline. 437

TABLE V
THE ABLATION STUDIES OF THE PROPOSED CADNET. THE AUC, ACC,
AP, AND EER ARE USED AS EVALUATION METRICS. THE BEST RESULTS
ARE HIGHLIGHTED IN BOLD.

Methods	AUC↑	ACC↑	AP↑	EER↓
Baseline	0.9831	0.9445	0.9960	0.0589
Baseline + CSC	0.9842	0.9455	0.9961	0.0547
Baseline + FDD	0.9841	0.9460	0.9962	0.0587
Baseline + CSC + FDD	0.9854	0.9479	0.9965	0.0576

V. CONCLUSION

In this work, the CADNet is proposed to enhance the 439 generalization ability of deepfake detection in the facial pay-440 ment system. In the proposed CADNet, a CSC module is 441 designed to focus on local features while emphasizing non-442 local facial regions. It can enhance the ability of the network 443 to detect general forgery traces. Meanwhile, an FDD module 444 is designed to decouple high-frequency and low-frequency 445 features for the exploration of frequency domain features. It 446 enables the network to focus on features that are overlooked 447 in the spatial domain. Experimental results demonstrate that 448 the proposed CADNet maintains strong generalization capa-449 bility. It enhances resistance to deepfake technology in finan-450 cial transactions while mitigating security threats posed by 451 such attacks. While CADNet excels at image-level detection, 452 its current single-frame analysis paradigm cannot inherently 453 identify temporal artifacts in video forgeries. Future work 454 will extend this framework to video deepfake detection by 455 analyzing inconsistencies between frames. 456

ACKNOWLEDGEMENTS

This work has been funded by the Shandong Province 458 Undergraduate Teaching Reform Project (No.Z2024184).

CONFLICT OF INTEREST STATEMENT 460

The authors have no conflict of interest related to this 461 publication. 462

DATA AVAILABILITY STATEMENT

Data sharing is not applicable to this article.

REFERENCES

- [1] Q. Li, M. Gao, G. Zhang, W. Zhai, J. Chen, and G. Jeon, "Towards 466 multimodal disinformation detection by vision-language knowledge interaction," Information Fusion, vol. 102, p. 102037, February 2024. 468
- [2] G. Zhang, M. Gao, Q. Li, W. Zhai, and G. Jeon, "Multi-modal generative deepfake detection via visual-language pretraining with gate fusion for 470 cognitive computation," Cognitive Computation, vol. 16, no. 6, pp. 471 2953-2966, 2024.
- [3] G. Zhang, M. Gao, Q. Li, S. Guo, and G. Jeon, "Detecting sequential 473 deepfake manipulation via spectral transformer with pyramid attention 474 in consumer iot," IEEE Transactions on Consumer Electronics, 2024. 475
- [4] M. S. Rana and A. H. Sung, "Deepfakestack: A deep ensemblebased learning technique for deepfake detection," in 2020 7th IEEE international conference on cyber security and cloud computing (CSCloud)/2020 6th IEEE international conference on edge computing and scalable cloud (EdgeCom). IEEE, 2020, pp. 70-75.
- [5] Z. Yan, Y. Zhang, Y. Fan, and B. Wu, "Ucf: Uncovering common features for generalizable deepfake detection," in 2023 IEEE/CVF International Conference on Computer Vision (ICCV), 2023, pp. 22355-22 366.

438

459

463

464

465

467

469

472

476

477

478

479

480

481

482

483

484

- [6] H. Liu, X. Li, W. Zhou, Y. Chen, Y. He, H. Xue, W. Zhang, and N. Yu,
 "Spatial-phase shallow learning: Rethinking face forgery detection in
 frequency domain," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 772–
 781.
- Y. Qian, G. Yin, L. Sheng, Z. Chen, and J. Shao, "Thinking in frequency: Face forgery detection by mining frequency-aware clues," in *Computer Vision – ECCV 2020*. Cham: Springer International Publishing, 2020, pp. 86–103.
- Y. Luo, Y. Zhang, J. Yan, and W. Liu, "Generalizing face forgery detection with high-frequency features," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 16317–16326.
- [9] H. Cheng, Y. Guo, T. Wang, L. Nie, and M. Kankanhalli, "Towards generalizable deepfake detection by primary region regularization," *arXiv preprint arXiv:2307.12534*, 2023.
- [10] G. Zhang, M. Gao, Q. Li, W. Zhai, G. Zou, and G. Jeon, "Disrupting deepfakes via union-saliency adversarial attack," *IEEE Transactions on Consumer Electronics*, vol. 70, no. 1, pp. 2018–2026, 2023.
- [11] G. Zhang, Q. Li, M. Gao, and G. Jeon, "Towards sequential deepfake
 detection using deep learning for privacy protection," *IEEE Consumer Electronics Magazine*, vol. 14, pp. 42–48, 2024.
- [12] S. Guo, Q. Li, M. Gao, G. Zhang, J. Pan, and G. Jeon, "Deep learning-based face forgery detection for facial payment systems," *IEEE Consumer Electronics Magazine*, vol. 14, no. 3, pp. 80–86, 2024.
- [13] S. Guo, Q. Li, M. Gao, G. Zhang, and G. Jeon, "Smart city security: Fake
 news detection in consumer electronics," *IEEE Consumer Electronics Magazine*, pp. 1–7, 2024.
- [14] Z. Yang, J. Liang, Y. Xu, X.-Y. Zhang, and R. He, "Masked relation learning for deepfake detection," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 1696–1708, 2023.
- [15] M.-Y. Tsai, H.-H. Cho, C.-M. Yu, Y.-C. Chang, and H.-C. Chao, "Effective adversarial examples identification of credit card transactions," *IEEE Intelligent Systems*, vol. 39, no. 4, pp. 50–59, 2024.
- [16] M. T. Usman, H. Khan, S. K. Singh, M. Y. Lee, and J. Koo, "Efficient deepfake detection via layer-frozen assisted dual attention network for consumer imaging devices," *IEEE Transactions on Consumer Electronics*, pp. 1–1, 2024.
- J. Tian, C. Yu, X. Wang, P. Chen, Z. Xiao, J. Dai, J. Han, and Y. Chai,
 "Real appearance modeling for more general deepfake detection," in *European Conference on Computer Vision*. Springer Nature Switzerland, 2025, pp. 402–419.
- M. Wu, J. Ma, R. Wang, S. Zhang, Z. Liang, B. Li, C. Lin, L. Fang, and
 L. Wang, "Traceevader: Making deepfakes more untraceable via evading
 the forgery model attribution," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 18, 2024, pp. 19965–19973.
- [19] S. Guo, Q. Li, M. Gao, X. Zhu, and I. Rida, "Generalizable deepfake detection via spatial kernel selection and halo attention network," *Image and Vision Computing*, vol. 160, p. 105582, July 2025.
- J. Cao, C. Ma, T. Yao, S. Chen, S. Ding, and X. Yang, "End-toend reconstruction-classification learning for face forgery detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 4113–4122.
- [21] Y. Ni, D. Meng, C. Yu, C. Quan, D. Ren, and Y. Zhao, "Core: Consistent representation learning for face forgery detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (CVPR) Workshops, June 2022, pp. 12–21.
- 542 [22] H. Dang, F. Liu, J. Stehouwer, X. Liu, and A. K. Jain, "On the
 543 detection of digital face manipulation," in *Proceedings of the IEEE/CVF*544 *Conference on Computer Vision and Pattern recognition*, 2020, pp.
 545 5781–5790.
- W. Lu, L. Liu, B. Zhang, J. Luo, X. Zhao, Y. Zhou, and J. Huang,
 "Detection of deepfake videos using long-distance attention," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [24] Y. Wang, K. Yu, C. Chen, X. Hu, and S. Peng, "Dynamic graph learning
 with content-guided spatial-frequency relation reasoning for deepfake
 detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7278–7287.
- [25] S. Das, M. S. Islam, and M. R. Amin, "Gca-net: utilizing gated context attention for improving image forgery localization and detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 81–90.
- F. Wang, Q. Chen, B. Jing, Y. Tang, Z. Song, and B. Wang, "Deepfake detection based on the adaptive fusion of spatial-frequency features," *International Journal of Intelligent Systems*, vol. 2024, no. 1, p. 7578036, 2024.

- [27] J. Gao, Z. Xia, G. L. Marcialis, C. Dang, J. Dai, and X. Feng, "Deepfake detection based on high-frequency enhancement network for highly compressed content," *Expert Systems with Applications*, vol. 249, p. 123732, September 2024.
- [28] J. Liu, J. Wang, P. Zhang, C. Wang, D. Xie, and S. Pu, "Multi-scale wavelet transformer for face forgery detection," in *Proceedings of the Asian Conference on Computer Vision*, 2022, pp. 1858–1874.
- [29] M. Wolter, F. Blanke, R. Heese, and J. Garcke, "Wavelet-packets for deepfake image analysis and detection," *Machine Learning*, vol. 111, no. 11, pp. 4295–4327, 2022.
- [30] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1501–1510.
- [31] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, "Faceforensics++: Learning to detect manipulated facial images," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 1–11.
- [32] Y. Li, X. Yang, P. Sun, H. Qi, and S. Lyu, "Celeb-df: A largescale challenging dataset for deepfake forensics," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 3207–3216.
- [33] B. Dolhansky, J. Bitton, B. Pflaum, J. Lu, R. Howes, M. Wang, and C. C. Ferrer, "The deepfake detection challenge (dfdc) dataset," *arXiv* preprint arXiv:2006.07397, 2020.
- [34] L. Li, J. Bao, H. Yang, D. Chen, and F. Wen, "Advancing high fidelity identity swapping for forgery detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 5074–5083.
- [35] Z. Yan, Y. Zhang, X. Yuan, S. Lyu, and B. Wu, "Deepfakebench: A comprehensive benchmark of deepfake detection," in *Advances in Neural Information Processing Systems*, vol. 36. Curran Associates, Inc., 2023, pp. 4534–4565.
- [36] D. Afchar, V. Nozick, J. Yamagishi, and I. Echizen, "Mesonet: a compact facial video forgery detection network," in 2018 IEEE International Workshop on Information Forensics and Security (WIFS), 2018, pp. 1–7.
- [37] S.-Y. Wang, O. Wang, R. Zhang, A. Owens, and A. A. Efros, "Cnngenerated images are surprisingly easy to spot... for now," in *Proceedings* of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 8695–8704.
- [38] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 09–15 Jun 2019, pp. 6105–6114.
- [39] H. H. Nguyen, J. Yamagishi, and I. Echizen, "Capsule-forensics: Using capsule networks to detect forged images and videos," in *ICASSP 2019* - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2019, pp. 2307–2311.
- [40] Y. Li, "Exposing deepfake videos by detecting face warping artif acts," *arXiv preprint arXiv:1811.00656*, 2018.
- [41] L. Li, J. Bao, T. Zhang, H. Yang, D. Chen, F. Wen, and B. Guo, "Face x-ray for more general face forgery detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 5001–5010.
- [42] D. Kinga, J. B. Adam *et al.*, "A method for stochastic optimization," in *International conference on learning representations (ICLR)*, vol. 5. San Diego, California; 2015, p. 6.
- [43] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, pp. 8024–8035, 2019.

561

562

563

564

565

617

618

619

620