

Defending Deepfake via Texture Feature Perturbation

Xiao Zhang¹, Changfang Chen² and Tianyi Wang^{3,*}

Abstract—The rapid development of Deepfake technology poses severe challenges to social trust and information security. While most existing detection methods primarily rely on passive analyses, due to unresolvable high-quality Deepfake contents, proactive defense has recently emerged by inserting invisible signals in advance of image editing. In this paper, we introduce a proactive Deepfake detection approach based on facial texture features. Since human eyes are more sensitive to perturbations in smooth regions, we invisibly insert perturbations within texture regions that have low perceptual saliency, applying localized perturbations to key texture regions while minimizing unwanted noise in non-textured areas. Our texture-guided perturbation framework first extracts preliminary texture features via Local Binary Patterns (LBP), and then introduces a dual-model attention strategy to generate and optimize texture perturbations. Experiments on CelebA-HQ and LFW datasets demonstrate the promising performance of our method in distorting Deepfake generation and producing obvious visual defects under multiple attack models, providing an efficient and scalable solution for proactive Deepfake detection.

I. INTRODUCTION

With the continuous improvement of generative algorithms, modern Deepfake technologies [1]–[3] can produce synthetic images that are indistinguishable to human eyes, with highly realistic facial features, voice timbre, and body movements. However, malicious usages of Deepfake have posed serious and urgent threats to the authenticity of media information [4].

Countermeasures for detection has been proposed since the first occurrence of Deepfake contents [5]. They generally fall into two categories: passive detection and proactive defense. Traditional passive detection [6]–[9] primarily focuses on identifying whether content is generated by Deepfake techniques. In contrast, proactive defense emphasizes preventing Deepfake models from generating convincing forged content by embedding imperceptible defense information or changing data features before the forged content is generated.

This study is grateful to the Key R&D Program of Shandong Province (Major Scientific and Technological Innovation Project): (Grant No.2023CXGC010113); The Taishan Scholars Program: (Grant NO.tspd20240814); Qilu University of Technology (Shandong Academy of Sciences) Project (Grant No.2024ZDZX11); Key Laboratory of Computing Power Network and Information Security, Ministry of Education, Qilu University of Technology (Shandong Academy of Sciences).

*Corresponding author: Tianyi Wang.

¹Xiao Zhang is with the Key Laboratory of Computing Power Network and Information Security, Ministry of Education, Qilu University of Technology (Shandong Academy of Sciences), Jinan, China 10431230097@stu.qlu.edu.cn

²Changfang Chen is with the Shandong Artificial Intelligence Institute, Qilu University of Technology (Shandong Academy of Sciences), Jinan, China chenbht012@qlu.edu.cn

³Tianyi Wang is with the School of Computing, National University of Singapore wangty@nus.edu.sg

Existing proactive defense methods are primarily divided into watermark-based [10]–[14] and perturbation-based [15]–[18] approaches. Watermarking embeds identifiable markers for copyright protection and traceability, however, the embedded information can be fragile under white-box attacks, making it susceptible to removal or overwriting by adversaries. In contrast, perturbation-based methods directly distort the generation results of Deepfake models by introducing subtle modifications that are imperceptible to the human eye. Existing perturbation defense methods that perform perturbation embedding in the spatial domain typically adopt a uniform noise injection strategy for the image. However, this indiscriminate perturbation pattern has two key flaws. First, since the human visual system is more sensitive to noise in flat areas than in edges or textured regions [19], the uniform injection strategy will cause noise to produce obvious visual artifacts in smooth areas, seriously affecting image quality. On the other hand, this method may result in insufficient perturbation intensity in key facial texture areas, making it difficult to effectively suppress the generative model’s ability to imitate these areas.

To address the aforementioned issues, this study proposes a perturbation framework guided by texture features. For an input image, we first preprocess it through bilateral filtering, extract basic texture features through local binary patterns in the texture extraction module, and generate attention feature maps using the gradient-weighted class activation mapping (Grad-CAM) technique. These two feature maps are then fused in the perturbation enhancement module, and a perturbed image is generated through multi-layer deformable convolution operations. Thereafter, the forged image and its corresponding attention feature map are derived. To optimize the perturbation effect, we construct a multi-objective loss function by comparing the feature differences between the original and the perturbed output. This effectively interferes with the forged generation process while ensuring visual quality. The contributions of this work are threefold:

- We propose a perturbation framework based on facial texture features. By accurately identifying and perturbing the texture features of key facial areas, the framework effectively destroys the imitation and generation of key textures in the Deepfake generation process.
- We design a dual-model attention strategy to guide the optimization of perturbations. This method integrates initial texture features with attention maps generated through local feature parsing for regional enhancement while leveraging global semantic modeling to optimize perturbation directions. By optimizing perturbations

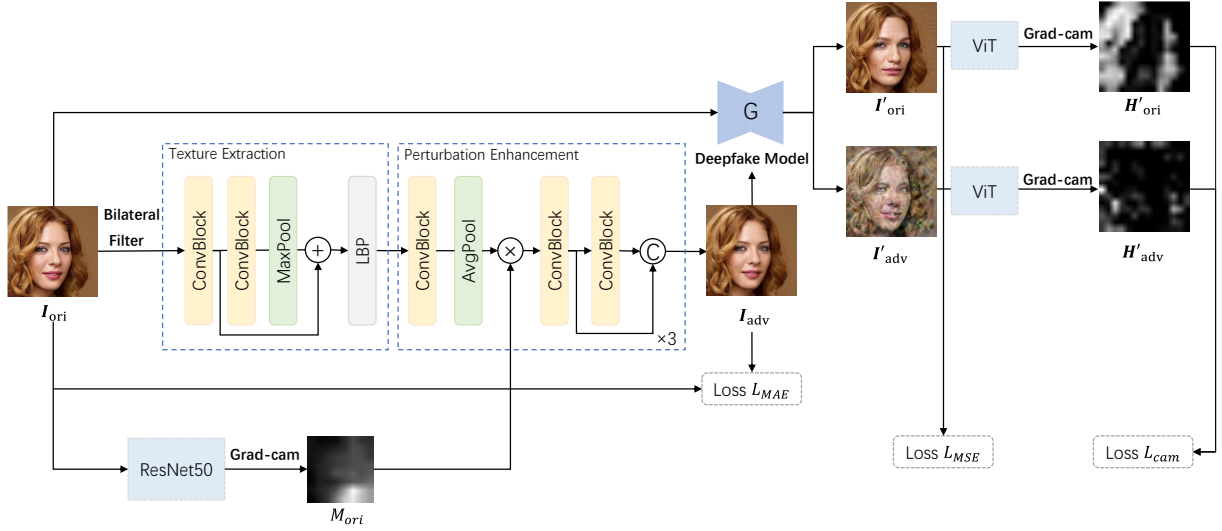


Fig. 1. Demonstration of the proposed framework. For the input image I_{ori} , the texture extraction module extracts its initial texture information and inputs it into the ResNet50 network, where Grad-CAM is applied to generate the corresponding feature attention map M_{ori} , the obtained M_{ori} and the initial texture information are enhanced in the perturbation enhancement module. The generated perturbed image exhibits obvious visual defects after passing through the Deepfake model.

across multiple spatial scales, we achieve a balance between defense effectiveness and visual quality.

- Extensive experiments demonstrate that the framework can effectively resist Deepfake manipulation models and produce perturbed results with significant visual defects.

II. RELATED WORK

A. Deepfake Generation

Deepfake attribute editing modifies specific facial attributes in images or videos while preserving the person's identity, achieving high-quality and natural results. StarGAN [20] employs a unified multi-domain architecture for cross-attribute image conversion without separate domain-pair training. AttentionGAN [21] enhances local detail editing by using spatial attention maps to isolate target regions, while AttGAN [22] disentangles attribute and identity features, enforcing target attributes via classification loss. HiSD [23] achieves fine-grained control through hierarchical style disentanglement, separating semantically relevant and irrelevant layers.

For expression reenactment, StarGAN-V2 [24] eliminates the need for explicit attribute labels by transferring styles from a source image to a target reference image, enabled by its mapping network and style encoder. HyperReenact [25] further refines this by leveraging StyleGAN2's [26] photorealistic generation, using a hypernetwork to manipulate latent spaces for precise identity preservation and pose retargeting, avoiding artifacts from external tools.

B. Deepfake Defense

Deepfake generation defense strategies can be divided into passive detection and proactive defense. Passive detection [27]–[31] typically determines the authenticity of an image by identifying anomalies in the spatial or frequency

domain. To overcome the shortcomings of passive detection, perturbation techniques use specific technical means in the Deepfake generation process. These methods intentionally introduce distortions and forgeries, effectively reducing the quality of synthetic images.

CMUA-Watermark [17] aims to generate universal adversarial perturbations that are effective across multiple models without relying on any specific Deepfake model. Anti-Forgery [32] generates adversarial perception-sensitive perturbations in the Lab color space, which exhibit strong robustness against various input transformation operations. TAFIM [15] optimizes a global perturbation pattern and a conditional generative model, where an attention network and a fusion network combine multiple model-specific perturbations to generate perturbations tailored to specific images. Meanwhile, TCA-GAN [33] introduces a transferable adversarial attack that effectively disrupts Deepfake models in black-box settings. Guan et al. [34] proposed an ensemble defense strategy called Hard Model Mining (HMM). This strategy does not optimize for good average performance during the perturbation optimization process. Instead, it improves the perturbation's ability to disrupt the most resilient models at each iteration step, thereby enabling it to disrupt all Deepfake models. To further enhance the generalization and robustness of perturbations, ComGAN [35] was proposed to learn the shared characteristics of compression across different platforms to generate robust adversarial perturbations.

III. METHODOLOGY

A. Method Overview

To satisfy the dual objectives of adversarial perturbation and visual quality, in this paper, we propose a perturbation method based on facial texture features. Unlike traditional uniform perturbation strategies, we embed perturbations

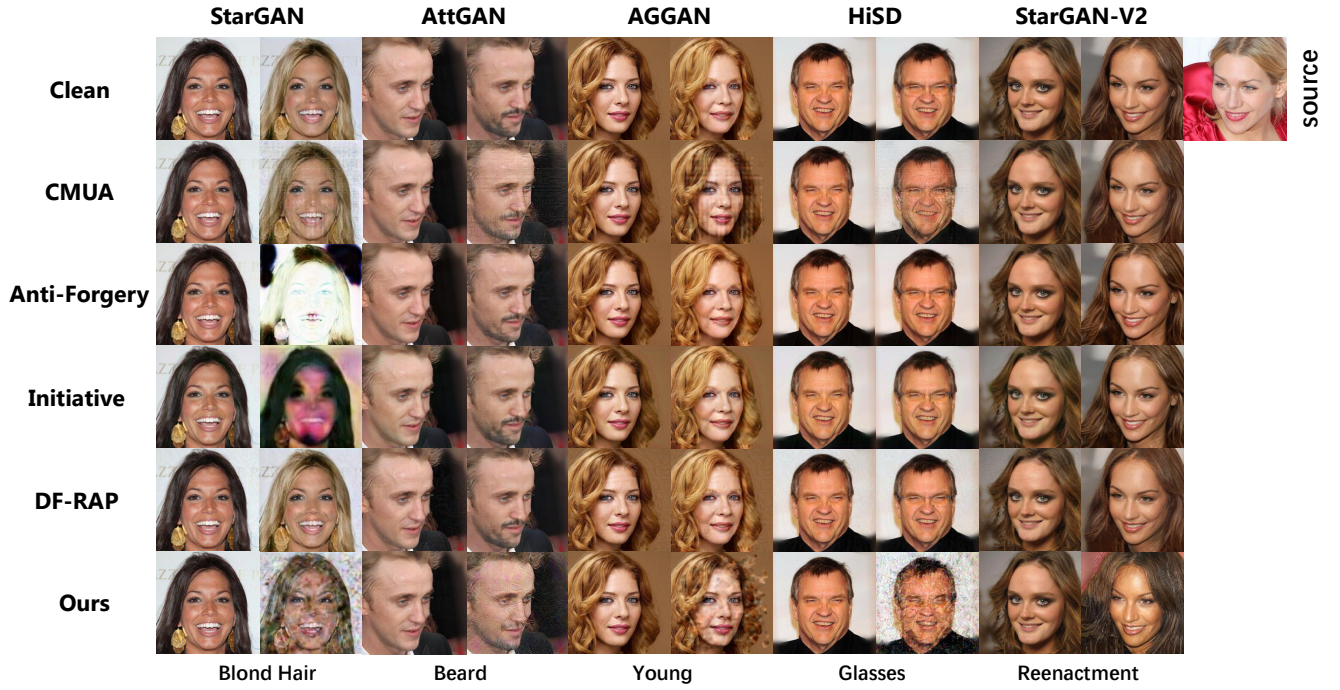


Fig. 2. Visual effects of Deepfake attribute editing and reenactment image manipulations on the different defense methods. The source image in the last column is dedicated to providing expression features for StarGAN-V2 when reenacting. In the remaining columns, all odd-numbered columns display perturbation images while even-numbered columns show the corresponding Deepfake results.

within the image’s textured regions. The naturally complex structures and high-frequency details in these regions provide an ideal context for concealing subtle perturbations or modifications. Furthermore, perturbations in texture regions disrupt the learning process of generative models on these critical features, causing distortions of the generated images and thereby degrading the overall generation quality.

Fig. 1 shows the design of the model. The input raw image I_{ori} is processed through two parallel branches. In the first branch, bilateral filtering [36] is applied to eliminate high-frequency noise interference, followed by a texture extraction module to obtain deep texture features. In the second branch, I_{ori} is fed directly into a ResNet50 network [37], where the Grad-CAM algorithm [38] is employed to generate the feature attention distribution M_{ori} . Subsequently, a perturbation enhancement module is designed to spatially enhance the initial texture features guided by the feature attention mechanism. This process ultimately produces visually imperceptible adversarial perturbations that can produce noticeable distortions after being tampered with by the Deepfake model.

B. Texture Extraction

We use Bilateral Filter combined with Local Binary Pattern (LBP) [39] for texture extraction. The core steps include grayscale conversion, bilateral filter preprocessing, and LBP texture quantization.

Grayscale Conversion. Texture feature analysis is typically based on luminance information. Therefore, we convert the original image $I_{\text{ori}} \in \mathbb{R}^{H \times W \times 3}$ to grayscale to eliminate color interference, using weighted average to compute the

luminance channel:

$$I_{\text{gray}}(x, y) = 0.299R(x, y) + 0.587G(x, y) + 0.116B(x, y), \quad (1)$$

where (x, y) represents the pixel coordinates, and the weight coefficients comply with the ITU-R BT.709 standard [40], ensuring that the converted grayscale image I_{gray} aligns with the human visual brightness perception characteristics.

Bilateral Filter. We apply a bilateral filter to the grayscale image in order to smooth noise while preserving texture edges. The filtered image I_{filter} suppresses noise in flat regions while enhancing high-frequency texture details. The bilateral filtering formula is denoted as follows:

$$I_{\text{filter}}(x, y) = \frac{1}{W_{\text{BF}}} \sum_{i, j \in \Omega} G_{\sigma_d}(i, j) \cdot G_{\sigma_r}(\Delta I) \cdot I_{\text{gray}}(x+i, y+j), \quad (2)$$

$$\Delta I = |I_{\text{gray}}(x, y) - I_{\text{gray}}(x+i, y+j)|, \quad (3)$$

where Ω represents the neighborhood window range, σ_d is the standard deviation of the spatial Gaussian kernel, and σ_r is the standard deviation of the range Gaussian kernel.

LBP Texture Quantization. The core of this process is to apply two ConvBlocks to the filtered image I_{filter} to extract multi-level structural features. After compressing the spatial information using max pooling, the LBP method is employed, with the formula as follows:

$$LBP_{P,R}(x_c, y_c) = \sum_{p=0}^{P-1} s(g_p - g_c) 2^p, \quad (4)$$

where P denotes the number of sampling points in the circular neighborhood, R denotes the radius of the neighborhood, g_c is the grayscale value of the central pixel, g_p is the grayscale value of the p -th sampling point, and x_c and y_c represent the coordinates of the central pixel in the image. When traversing the entire image, each pixel becomes the central pixel (x_c, y_c) once, and its LBP value is calculated using this formula. $s(x)$ is the sign function, defined as follows:

$$s(x) = \begin{cases} 1 & \text{if } x \geq 0, \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

C. Perturbation Enhancement

In order to ensure that the perturbation focuses on the texture regions while reducing the unwanted noise in the smooth area, we introduced an attention guidance mechanism based on Grad-CAM, and enhanced the initial texture area through a ConvBlock. A single ConvBlock consists of a convolutional neural network (CNN) layer, a batch normalization layer, and a ReLU activation function. The feature map output by the convolution block is spatially compressed using average pooling. At the same time, we use the pre-trained classification model to forward propagate the input image I_{ori} , generate the feature attention map M_{ori} through reverse gradient calculation, and perform channel-by-channel weighted fusion with the multi-scale texture features after spatial compression to ensure that the perturbation is aligned with the model's attention area. This fusion strategy ensures that the perturbation enhancement is always spatially aligned with the key areas of the model's decision, thereby significantly improving the targeting of adversarial attacks.

D. Objective Functions

In this study, we aim to effectively prevent the generation of Deepfake while forming perturbed images with high visual quality. To maintain perceptual fidelity during perturbation, we formulate the optimization problem using an L_1 distance metric between corresponding pixels, which is defined as:

$$\mathcal{L}_{\text{MAE}} = \|I_{\text{adv}} - I_{\text{ori}}\|_1, \quad (6)$$

where I_{ori} represents the original image, and I_{adv} represents the perturbed image.

The L_2 norm is adopted as a constraint to regulate the differences between the unperturbed and perturbed images, both in their original and generated forms. The loss function is defined as

$$\mathcal{L}_{\text{MSE}} = -\|I'_{\text{ori}} - I'_{\text{adv}}\|_2^2, \quad (7)$$

where I'_{ori} and I'_{adv} represent the generated results of I_{ori} and I_{adv} after the generative model is applied.

Adversarial Attention Loss function dynamically focuses on critical regions through a difference region mask, maximizing the attention discrepancy in these key areas to optimize perturbation generation. This ensures that the generated fake images exhibit significant differences in the attention regions of the Vision Transformer (ViT) [41] following

$$\mathcal{L}_{\text{cam}} = -\log \left(\frac{\sum M_{\text{diff}} \cdot |H'_{\text{ori}} - H'_{\text{adv}}|}{\sum M_{\text{diff}}} \right), \quad (8)$$

and

$$M_{\text{diff}} = \mathbb{I}(|H'_{\text{ori}} - H'_{\text{adv}}| > T), \quad (9)$$

where H'_{ori} and H'_{adv} are the Grad-CAM heatmaps generated by the ViT network, corresponding to the generated image I'_{ori} and the perturbed image I'_{adv} , respectively. M_{diff} is a binary mask that highlights regions where the absolute difference between H'_{ori} and H'_{adv} exceeds a threshold T .

The total loss is computed as:

$$\mathcal{L}_{\text{total}} = \lambda_1 \mathcal{L}_{\text{MAE}} + \lambda_2 \mathcal{L}_{\text{MSE}} + \lambda_3 \mathcal{L}_{\text{cam}}, \quad (10)$$

where λ_1 , λ_2 , and λ_3 are the weights that determine the influence of each loss component.

TABLE I

QUANTITATIVE VISUAL QUALITY EVALUATION OF PERTURBED IMAGES ON THE CELEBA-HQ DATASET. BEST PERFORMANCE MARKED IN **BOLD**.

Model	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
CMUA [17]	38.6395	0.9504	0.0333
Anti-Forgery [32]	38.0704	0.9529	0.0281
Initiative [42]	35.8042	0.9114	0.0866
DF-RAP [35]	38.8466	0.9349	0.0511
Ours	39.9355	0.9617	0.0251

IV. EXPERIMENTS

A. Implementation Details and Metrics

Deepfake Models. We validated our framework against leading adversarial models, including StarGAN [20], AttGAN [22], AGGAN [21], and HiSD [23] for attribute editing and StarGAN-V2 [24] for expression reenactment. StarGAN and AGGAN were trained for multi-attribute facial editing, covering five key attributes: black hair, blond hair, brown hair, gender, and age. AttGAN supports a broader range of 14 attributes, including skin color, bangs, beard, etc., while HiSD specializes in glasses manipulation.

Datasets. We conducted experiments using the CelebA-HQ [43] dataset, which consists of 30,000 high-resolution (1024 \times 1024) image samples and features 6,217 unique identities, covering diverse facial attributes. The dataset was partitioned into training, validation, and testing sets following the official split. For cross-dataset evaluation, we utilized the LFW [44] dataset, a widely adopted benchmark for unconstrained face verification and recognition. The LFW dataset contains 5,749 distinct identities, making it suitable for assessing generalization performance. In our experiments, all facial images are resized to 256 \times 256.

Parameters. For the bilateral filter, $\Omega = 31$ represents the neighborhood window range, $\sigma_d = 75$ is the standard deviation of the spatial Gaussian kernel, and $\sigma_i = 15$ is the standard deviation of the range Gaussian kernel. We selected $P = 8$ sampling points and a neighborhood radius of $R = 1$ to compute the LBP feature map. The threshold T in Eqn.(9) is set to 0.3. For the coefficients in Eqn.(10), we set $\lambda_1 = 1.0$, $\lambda_2 = 0.04$, and $\lambda_3 = 0.1$ for total loss.

Evaluation Metrics. We evaluated the visual quality of perturbed images using three metrics, peak signal-to-noise

TABLE II
QUALITATIVE EVALUATION OF THE MEAN L_2 NORM DISTANCE AND DSR BETWEEN I'_{ORI} AND I'_{ADV} ON THE CELEBA-HQ DATASET. THE BEST RESULT IS MARKED IN **BOLD**.

Models	CMUA [17]		Anti-Forgery [32]		Initiative [42]		DF-RAP [35]		Ours	
	$L_2\uparrow$	DSR \uparrow	$L_2\uparrow$	DSR \uparrow	$L_2\uparrow$	DSR \uparrow	$L_2\uparrow$	DSR \uparrow	$L_2\uparrow$	DSR \uparrow
StarGAN [20]	0.045	35.46%	0.484	100%	0.152	98.80%	0.006	0.23%	0.340	100%
AttGAN [22]	0.041	42.82%	0.017	0.07%	0.010	0.0%	0.002	0.0%	0.046	51.22%
AGGAN [21]	0.041	48.01%	0.004	0.02%	0.006	0.08%	0.003	0.02%	0.161	100%
HiSD [23]	0.052	54.77%	0.004	0.04%	0.016	0.47%	0.007	0.10%	0.079	64.79%
StarGAN-V2 [24]	0.029	14.83%	0.040	23.23%	0.038	21.28%	0.027	12.28%	0.097	93.06%
Average	0.042	39.18%	0.110	24.67%	0.044	24.13%	0.009	2.53%	0.145	81.81%

ratio (PSNR) to measure pixel-level fidelity, structural similarity index (SSIM) to assess structural preservation, and learned perceptual image patch similarity (LPIPS) to quantify perceptual differences aligned with human visual perception. To evaluate the effectiveness of the proposed defense method, we computed the visual distance between the original tampered image and perturbation-protected tampered image using the L_2 norm. Additionally, we utilized the Defense Success Rate (DSR), which records the percentage of disrupted images exhibiting distortions with $L_2 \geq 0.05$ as a metric to measure the efficacy of the defense.

B. Visual Quality

Table I presents a comparative evaluation of five defense methods across three image quality metrics: PSNR, SSIM, and LPIPS. Our method consistently achieves superior performance across all metrics. Specifically, it attains the highest PSNR value of 39.9355, being the only method among the compared approaches to surpass 39 dB. For structural similarity, our method achieves an SSIM score of 0.9617, outperforming Anti-Forgery by 0.0088. In terms of perceptual similarity, we further reduce the LPIPS score by 10.7% compared to the second-best method, indicating higher visual fidelity and better alignment with human perception.

Fig. 2 visualizes the perturbed images generated by each method and the corresponding Deepfake-forged images. Specifically, in order to intuitively compare the visual quality of the perturbed images and the defense effect of the Deepfake generated images, we showed the original input images and the normal Deepfake results in the first row. It is worth noting that we add a source image in the last column, which is used to provide expression features in StarGAN-V2. Every two remaining columns show the perturbed images generated by different perturbation methods and the defense effects of the corresponding perturbed images after attribute editing or expression reenactment. It can be observed that the perturbed images generated by the Initiative method present an unnatural green effect, the perturbed images from DR-RAP display unnatural artifacts that are visually noticeable, and our method introduces obvious distortions while maintaining excellent visual quality.

C. Comparison with SOTA

In Table II, we compared the mean L_2 norm distance and DSR performance with four popular perturbation methods, namely, CMUA [17], Anti-Forgery [32], Initiative [42], and DF-RAP [35]. As demonstrated in Table II, our method achieves superior performance compared to other state-of-the-art models across multiple metrics. Specifically, when defending against StarGAN attacks, both Anti-Forgery and Initiative exhibit strong defensive capabilities. Although the L_2 distance of our perturbation is slightly lower than Anti-Forgery’s, it still achieves a 100% defense success rate, matching the best existing perturbation methods. When fighting against AttGAN, our method fails to cause serious distortions in the results, but it produces perceptible manipulation traces in the visualizations. In terms of the average attack effect of the five models, our method surpasses the second-best method by 42.63% in defense success rate.

We used the attributes of blond hair, beard, young, and glasses to visualize the adversarial effects of StarGAN, AttGAN, AGGAN, and HiSD. Our perturbations induce visible distortions in the generated images across all Deepfake models, disrupting their ability to produce natural-looking results. For StarGAN-V2, it transfers facial expressions from a source image to a target image while preserving the target’s hairstyle and background. In our experimental design, all perturbations are added to the target image, so the generated results of StarGAN-V2 reenact the source expressions on the target images. The results show that none of the four defense methods we compared can resist the manipulation of StarGAN-V2, our approach uniquely compromised its generation process. Although our perturbations do not introduce noticeable artifacts in facial regions, they successfully cause severe distortions in hairstyles and background elements.

In order to further prove the effectiveness of our method on unseen datasets, we used LFW for cross-dataset verification. Since the LFW dataset does not include attribute annotations for the original images, we omitted the Anti-Forgery approach in our experiments. As shown in Table IV, our method effectively preserves the visual quality and texture structure of the original protected image, and can still show the best visual quality on unseen datasets. As shown in Table III,

TABLE III

QUALITATIVE EVALUATION OF THE MEAN L_2 NORM DISTANCE AND DSR BETWEEN I'_{ORI} AND I'_{ADV} ON THE LFW DATASET. THE BEST RESULT IS MARKED IN **BOLD**.

Models	CMUA [17]		Initiative [42]		DF-RAP [35]		Ours	
	$L_2\uparrow$	DSR \uparrow	$L_2\uparrow$	DSR \uparrow	$L_2\uparrow$	DSR \uparrow	$L_2\uparrow$	DSR \uparrow
StarGAN [20]	0.055	57.52%	0.035	29.50%	0.004	0.01%	0.177	100%
AttGAN [22]	0.038	12.57%	0.004	0.02%	0.009	0.11%	0.042	50.57%
AGGAN [21]	0.048	41.50%	0.043	40.84%	0.002	0.0%	0.072	90.75%
HiSD [23]	0.045	39.09%	0.003	0.0%	0.001	0.0%	0.069	60.32%
StarGAN-V2 [24]	0.056	47.55%	0.066	61.92%	0.030	14.69%	0.143	99.72%
Average	0.048	39.65%	0.030	26.46%	0.009	2.96%	0.101	80.27%

TABLE IV

QUANTITATIVE VISUAL QUALITY EVALUATION OF PERTURBED IMAGES ON THE LFW DATASET. THE BEST RESULT IS MARKED IN **BOLD**.

Model	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
CMUA [17]	38.9998	0.9497	0.0581
Initiative [42]	36.8727	0.8775	0.1325
DF-RAP [35]	37.1954	0.9179	0.1241
Ours	39.9101	0.9545	0.0195

CMUA shows a certain defense ability when fighting against five models, but neither the single defense nor the average defense effects can outperform our method. This result illustrates the generalization ability and robustness of our method, proving that it has practical application potential in the field of perturbation defense.

D. Ablation Study

We conducted separate experiments using ResNet50 and ViT to generate decision key maps for guiding perturbation generation. This single-model approach allows us to verify the effectiveness of each model's key components in perturbation generation, and examine how the Grad-CAM-guided attention mechanism operates in isolation. We first tested ResNet50 independently, then repeated the experiments using ViT alone.

Following the same evaluation process as the previous experiment, the results of each training are exhibited in Table V. In terms of visual quality, ResNet50-only maintains similar visual effects to our model, but ViT-only lacks the local attention guidance of ResNet50. The LBP texture feature cannot accurately locate the high-frequency areas that the model is sensitive to. The perturbation enhancement becomes a global uniform injection, resulting in the smooth area also being covered by the perturbation, affecting the visual quality. In terms of adversarial effects, the average defense success rate of ResNet50-only decreased by approximately 8% due to the lack of loss calculation of the difference area mask. While ResNet50-only achieves marginally higher visual quality scores, the improvement is trivial and visually undetectable. In contrast, our method achieves superior L_2 and DSR performance, establishing the optimal balance

between fidelity and defense performance.

TABLE V

THE VISUAL QUALITY AND DEFENSE EFFECTIVENESS UNDER DIFFERENT SETTINGS.

Models	ResNet50-only		ViT-only		Ours	
	$L_2\uparrow$	DSR \uparrow	$L_2\uparrow$	DSR \uparrow	$L_2\uparrow$	DSR \uparrow
StarGAN [20]	0.318	100%	0.405	100%	0.340	100%
AttGAN [22]	0.022	32.27%	0.042	50.35%	0.046	51.22%
AGGAN [21]	0.101	98.82%	0.017	100%	0.161	100%
HiSD [23]	0.052	56.14%	0.072	60.11%	0.079	64.79%
StarGAN-V2 [24]	0.082	81.25%	0.090	89.34%	0.097	93.06%
Average	0.115	73.70%	0.125	79.96%	0.145	81.81%
PSNR \uparrow	39.9547		37.5541		39.9355	
SSIM \uparrow	0.9688		0.9577		0.9617	
LPIPS \downarrow	0.0279		0.0446		0.0251	

V. CONCLUSIONS

In this paper, we propose a facial texture-aware perturbation generation framework with dual-branch collaborative optimization, which proactively defends against malicious Deepfake manipulations. After initially extracting texture features using LBP, we identify key image areas via a perturbation enhancement module integrated with Grad-CAM, and guide perturbation generation under local detail enhancement and global semantic constraints. This method achieves a balance between adversarial resistance and visual fidelity. Experiments demonstrate that our framework is effective in disrupting a variety of Deepfake models. The generated forged images exhibit significant distortions, while the original images retain imperceptible perturbations that do not affect human perception, demonstrating the robustness and versatility of our approach across different scenarios.

REFERENCES

- [1] R. Chen, X. Chen, B. Ni, and Y. Ge, "Simswap: An efficient framework for high fidelity face swapping," in *Proceedings of the 28th ACM international conference on multimedia*, 2020, pp. 2003–2011.
- [2] X. Ren, X. Chen, P. Yao, H.-Y. Shum, and B. Wang, "Reinforced disentanglement for face swapping without skip connection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 20 665–20 675.

- [3] J. Zhou, X. Jia, Q. Li, L. Shen, and J. Duan, "Uniface: Unified cross-entropy loss for deep face recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 20 730–20 739.
- [4] F. Juefei-Xu, R. Wang, Y. Huang, Q. Guo, L. Ma, and Y. Liu, "Countering malicious deepfakes: Survey, battleground, and horizon," *International journal of computer vision*, vol. 130, no. 7, pp. 1678–1734, 2022.
- [5] T. Wang, X. Liao, K. P. Chow, X. Lin, and Y. Wang, "Deepfake detection: A comprehensive survey from the reliability perspective," *ACM Computing Surveys*, vol. 57, no. 3, pp. 1–35, 2024.
- [6] Z. Guo, G. Yang, J. Chen, and X. Sun, "Fake face detection via adaptive manipulation traces extraction network," *Computer Vision and Image Understanding*, vol. 204, p. 103170, 2021.
- [7] T. Wang, H. Cheng, K. P. Chow, and L. Nie, "Deep convolutional pooling transformer for deepfake detection," *ACM transactions on multimedia computing, communications and applications*, vol. 19, no. 6, pp. 1–20, 2023.
- [8] T. Wang and K. P. Chow, "Noise based deepfake detection via multi-head relative-interaction," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 12, 2023, pp. 14 548–14 556.
- [9] S. Yang, J. Wang, Y. Sun, and J. Tang, "Multi-level features global consistency for human facial deepfake detection," *Journal of Image and Graphics*, vol. 27, no. 09, pp. 2708–2720, 2022.
- [10] N. Beuve, W. Hamidouche, and O. Déforges, "Waterlo: Protect images from deepfakes using localized semi-fragile watermark," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 393–402.
- [11] R. Wang, F. Juefei-Xu, M. Luo, Y. Liu, and L. Wang, "Faketagger: Robust safeguards against deepfake dissemination via provenance tracking," in *Proceedings of the 29th ACM international conference on multimedia*, 2021, pp. 3546–3555.
- [12] T. Wang, M. Huang, H. Cheng, B. Ma, and Y. Wang, "Robust identity perceptual watermark against deepfake face swapping," 2024.
- [13] T. Wang, M. Huang, H. Cheng, X. Zhang, and Z. Shen, "Lampmark: Proactive deepfake detection via training-free landmark perceptual watermarks," in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024, pp. 10515–10524.
- [14] T. Wang, H. Cheng, M.-H. Liu, and M. Kankanhalli, "Fractalforensics: Proactive deepfake detection and localization via fractal watermarks," 2025.
- [15] S. Aneja, L. Markhasin, and M. Nießner, "Tafim: Targeted adversarial attacks against facial image manipulations," in *European Conference on Computer Vision*. Springer, 2022, pp. 58–75.
- [16] Z. He, W. Wang, W. Guan, J. Dong, and T. Tan, "Defeating deepfakes via adversarial visual reconstruction," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 2464–2472.
- [17] H. Huang, Y. Wang, Z. Chen, Y. Zhang, Y. Li, Z. Tang, W. Chu, J. Chen, W. Lin, and K.-K. Ma, "Cmua-watermark: A cross-model universal adversarial watermark for combating deepfakes," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 1, 2022, pp. 989–997.
- [18] T. Wang, H. Cheng, X. Zhang, and Y. Wang, "Nullswap: Proactive identity cloaking against deepfake face swapping," 2025.
- [19] C. Luo, Q. Lin, W. Xie, B. Wu, J. Xie, and L. Shen, "Frequency-driven imperceptible adversarial attack on semantic similarity," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 15 315–15 324.
- [20] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo, "Stargan: Unified generative adversarial networks for multi-domain image-to-image translation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8789–8797.
- [21] H. Tang, D. Xu, N. Sebe, and Y. Yan, "Attention-guided generative adversarial networks for unsupervised image-to-image translation," in *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2019, pp. 1–8.
- [22] Z. He, W. Zuo, M. Kan, S. Shan, and X. Chen, "Attgan: Facial attribute editing by only changing what you want," *IEEE transactions on image processing*, vol. 28, no. 11, pp. 5464–5478, 2019.
- [23] X. Li, S. Zhang, J. Hu, L. Cao, X. Hong, X. Mao, F. Huang, Y. Wu, and R. Ji, "Image-to-image translation via hierarchical style disentanglement," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 8639–8648.
- [24] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha, "Stargan v2: Diverse image synthesis for multiple domains," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 8188–8197.
- [25] S. Bounareli, C. Tzelepis, V. Argyriou, I. Patras, and G. Tzimiropoulos, "Hyperreenact: one-shot reenactment via jointly learning to refine and retarget faces," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 7149–7159.
- [26] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, "Analyzing and improving the image quality of stylegan," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 8110–8119.
- [27] J. Li, H. Xie, J. Li, Z. Wang, and Y. Zhang, "Frequency-aware discriminative feature learning supervised by single-center loss for face forgery detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 6458–6467.
- [28] 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.
- [29] K. Shiohara and T. Yamasaki, "Detecting deepfakes with self-blended images," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 18 720–18 729.
- [30] Z. Xu, X. Zhang, R. Li, Z. Tang, Q. Huang, and J. Zhang, "Fakeshield: Explainable image forgery detection and localization via multi-modal large language models," *arXiv preprint arXiv:2410.02761*, 2024.
- [31] H. Zhao, W. Zhou, D. Chen, T. Wei, W. Zhang, and N. Yu, "Multi-attentional deepfake detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 2185–2194.
- [32] R. Wang, Z. Huang, Z. Chen, L. Liu, J. Chen, and L. Wang, "Anti-forgery: Towards a stealthy and robust deepfake disruption attack via adversarial perceptual-aware perturbations," *arXiv preprint arXiv:2206.00477*, 2022.
- [33] J. Dong, Y. Wang, J. Lai, and X. Xie, "Restricted black-box adversarial attack against deepfake face swapping," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 2596–2608, 2023.
- [34] W. Guan, Z. He, W. Wang, J. Dong, and B. Peng, "Defending against deepfakes with ensemble adversarial perturbation," in *2022 26th International Conference on Pattern Recognition (ICPR)*. IEEE, 2022, pp. 1952–1958.
- [35] Z. Qu, Z. Xi, W. Lu, X. Luo, Q. Wang, and B. Li, "Df-rap: A robust adversarial perturbation for defending against deepfakes in real-world social network scenarios," *IEEE Transactions on Information Forensics and Security*, 2024.
- [36] C. Tomasi and R. Manduchi, "Bilateral filtering for gray and color images," in *Sixth international conference on computer vision (IEEE Cat. No. 98CH36271)*. IEEE, 1998, pp. 839–846.
- [37] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [38] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [39] T. Ojala, M. Pietikainen, and D. Harwood, "Performance evaluation of texture measures with classification based on kullback discrimination of distributions," in *Proceedings of 12th international conference on pattern recognition*, vol. 1. IEEE, 1994, pp. 582–585.
- [40] I.-R. R. BT, "Parameter values for the hdtv standards for production and international programme exchange," *International Telecommunication Union, Recommendation*, 2002.
- [41] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [42] Q. Huang, J. Zhang, W. Zhou, W. Zhang, and N. Yu, "Initiative defense against facial manipulation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, 2021, pp. 1619–1627.
- [43] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," *arXiv preprint arXiv:1710.10196*, 2017.
- [44] G. Huang, M. Mattar, H. Lee, and E. Learned-Miller, "Learning to align from scratch," *Advances in neural information processing systems*, vol. 25, 2012.