ADAPTIVE GRAY: REDUCING COLOR DEPENDENCY TO IMPROVE GENERALIZATION IN DEEPFAKE DETECTION

Anonymous authorsPaper under double-blind review

ABSTRACT

Deepfake technology, powered by advanced generative models like GANs and diffusion models, has raised serious ethical and security concerns due to its potential for misuse in creating realistic yet deceptive content. These generative models are becoming increasingly sophisticated, making it harder for humans to distinguish real images from generated ones. This highlights the need for reliable machinebased detection. However, current detection methods face significant challenges in generalization, particularly when dealing with different generative models (crossgenerator) and diverse image scenarios (cross-dataset), such as faces, landscapes, and objects, limiting their applicability across various contexts. To address this challenge, we identified that the color dependency can often be unnecessary and may even impede deepfake detection performance. Building on this insight, we introduce Adaptive Gray (AG), a novel approach designed to improve classifier generalization by compressing the RGB channels of images. Our experiments on the large-scale GenImage dataset demonstrate that Adaptive Gray achieves the highest improvement of 19.9% in average ACC, 22.0% in AP, and 20.1% in TPR (at FPR=5%), consistently outperforming state-of-the-art classifiers. Meanwhile, inference efficiency improved by at most 1×10^4 times.

1 Introduction

The rapid proliferation of hyper-realistic images generated by advanced models like diffusion models Dhariwal & Nichol (2021); Ho et al. (2020); Nichol & Dhariwal (2021); Rombach et al. (2021); Mid (2022) poses growing societal threats, as humans increasingly struggle to distinguish synthetic content from genuine material Frank et al. (2024). Regulatory responses are emerging accordingly; for example, in September 2025 New South Wales (Australia) amended the *Crimes Act 1900* to criminalise the creation and sharing of sexually explicit deepfakes depicting real persons, with penalties up to three years' imprisonment (NSW, 2025). This highlights an urgent need for robust automated deepfake detection.

Despite significant research, current State-of-the-Art (SOTA) deepfake detection methods often fail to generalize effectively across unseen generative models (cross-generator) or diverse image domains (cross-dataset, e.g., faces vs. landscapes) Sha et al. (2023); Tan et al. (2024); Wang et al. (2023). Our in-depth analysis reveals that while some SOTA approaches show promise in constrained settings, their performance in real-world, dynamic scenarios remains unsatisfactory. This critical limitation stems from their failure to capture intrinsic, robust distinctions between real and synthetic images.

Our Work. We identify a crucial and often overlooked bottleneck hindering deepfake detection generalization: the detrimental impact of an over-reliance on color dependency. Our systematic empirical analysis reveals that conventional RGB representations, despite their information richness, introduce excessive color-dependent redundancy and noise. This noise can mislead classifiers, causing them to learn spurious correlations rather than robust, intrinsic generative artifacts. Consequently, models trained on full-color images struggle to adapt when color statistics shift across diverse datasets or novel generative models. Our key insight is that intelligently collapsing these

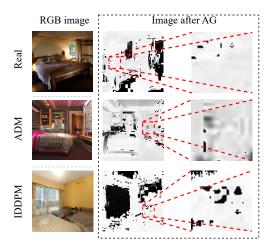


Figure 1: Qualitative comparison of artifacts in generated images before and after AG processing. The top row shows original RGB images, while the bottom row displays the same images after our proposed Adaptive Gray (AG) transformation. Notably, AG processing significantly enhances the visibility of texture-based artifacts inherent in generated images (e.g., ADM Dhariwal & Nichol (2021), IDDPM Nichol & Dhariwal (2021)) while suppressing color-dependent noise. This visual transformation makes distinguishing features more apparent and substantially improves classifier separation ability, particularly across diverse generative models.

color channels based on their intrinsic statistical differences can amplify hidden generative artifacts and reveal more general distinguishing features.

Inspired by this pivotal insight, we propose *Adaptive Gray* (AG), a novel, lightweight, and interpretable method. AG is designed to mitigate problematic color dependency and enhance generalization by intelligently transforming RGB images into a specialized grayscale representation. Unlike fixed grayscale conversions, AG learns optimal, data-driven coefficients to linearly combine RGB channels through an adaptive training process. This allows AG to effectively suppress irrelevant color details and accentuate subtle, texture-based artifacts unique to synthetic images, as visually demonstrated in Figure 1.

We conduct extensive experiments demonstrating AG's superior effectiveness and efficiency. Evaluated on challenging benchmark datasets like GenImage Zhu et al. (2024), AG consistently outperforms SOTA classifiers, achieving remarkable average improvements of 19.9% ACC, 22.0% AP, and 20.1% TPR@FPR=5%. Furthermore, AG significantly boosts inference efficiency by over 1×10^4 times. Our work not only highlights a fundamental challenge in deepfake detection but also offers a powerful, generalizable solution validated by substantial empirical gains.

Contributions. Our main contributions are summarized as follows:

- We uncover and systematically analyze a critical, overlooked factor affecting deepfake detection generalization: the detrimental impact of color dependency.
- Inspired by this insight, we propose *Adaptive Gray (AG)*, a novel, lightweight, and interpretable method designed to mitigate color dependency by learning optimal grayscale transformations.
- We conduct extensive experiments demonstrating AG's superior effectiveness, robustness, generalization, and remarkable inference efficiency across challenging datasets, reinforcing our initial hypothesis.

2 RELATED WORKS

The detection of generated images has been widely explored in recent years McCloskey & Albright (2018; 2019); Guo et al. (2018). As GANs and diffusion-based models have advanced, distinguish-

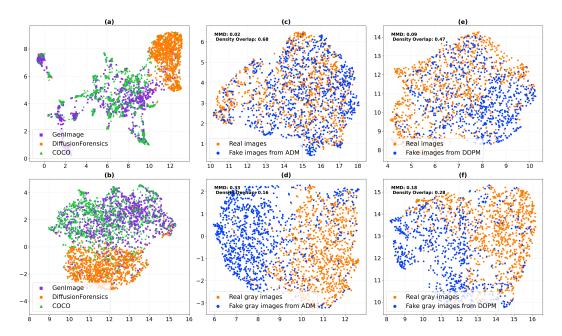


Figure 2: The distribution visualization before and after gray processing. (a) and (b) show dataset distribution changes before and after gray compression. (c) and (d) evaluate real/fake detection using ADM Dhariwal & Nichol (2021)-generated images, and (e) and (f) use DDPM Ho et al. (2020)-generated images—both comparing results before and after gray processing.

ing synthetic images from real ones has become increasingly difficult Karras et al. (2019); Brock (2018); Sauer et al. (2021).

Most existing methods focus on spatial artifacts introduced by generative models Yu et al. (2019); Marra et al. (2019); Ricker et al. (2024); Wang et al. (2020); Sarkar et al. (2024); Wang et al. (2023); Sha et al. (2023). For example, Wang et al. (2020) utilized deep neural networks to learn the distinguishing artifacts in generated images. However, their approach heavily relied on large amounts of GAN-generated images from diverse datasets. Sarkar Sarkar et al. (2024) conducted a noteworthy study focusing on the coherence of lines and shadows in images, providing insights into how physical inconsistencies could reveal generated content. DIRE Wang et al. (2023) used reconstruction residuals for detection, but its reliance on expensive image reconstruction made it impractical for large-scale detection tasks.

For fake face detection, researchers focus on identifying some representative artifacts based on the distinct features of facial images Rossler et al. (2019); Haliassos et al. (2021); Wang & Deng (2021). For example, Haliassos Haliassos et al. (2021) directed their work towards specific facial regions, such as the eyes and mouth, enhancing generated face detection by focusing on these distinctive areas. Wang Wang & Deng (2021) improved fake face detection by employing an attention-based data augmentation method to guide the detector to explore representative facial regions.

Although these studies in deepfake detection have yielded promising results when training and testing on fake images from the same distribution, they still struggle to generalize detectors to fake images from different distributions, whether generated by different models or sourced from different datasets. For example, Jeong Jeong et al. (2022) strengthened model performance by generating unique fingerprints for classification. Chen Chen et al. (2022) adopted adversarial training methods, enhancing model generalization and detection accuracy. Tan Tan et al. (2024) exploited upsampling artifacts through Neighboring Pixel Relationships (NPR), but their method struggled with datasets featuring diverse image distributions, as shown by our own experiments. While these approaches attempt to improve the generalization of deepfake detection for cross-model generalization, none of them focus on the challenge of cross-dataset generalization in deepfake detection.

Unlike previous methods, our method uniquely focuses on reducing color dependency in fake image detection, which allows us to avoid the need for computing new data representation features during

inference, as required by Wang Wang et al. (2023) and Tan Tan et al. (2024), thereby improving efficiency in the inference. Additionally, this streamlined approach significantly enhances the model's generalization capability (explained in Section 3.1).

3 METHODOLOGY

In this section, we introduce $Adaptive\ Gray\ (AG)$, a novel approach designed to enhance the generalization capabilities of deepfake detection models. Our methodology is fundamentally driven by a critical empirical observation: the often-detrimental impact of color dependency on model performance across diverse generative models and datasets. We begin by detailing the systematic analysis that led to this core insight, followed by the technical design of AG, which is engineered to effectively mitigate this issue through an adaptive, data-driven grayscale transformation process. Finally, we elaborate on the co-adaptive training strategy that enables AG to learn optimal image representations for robust deepfake detection.

3.1 EMPIRICAL EVIDENCE FOR COLOR DEPENDENCY

Our work is premised on the hypothesis that an unnecessary reliance on color information within deepfake detection models significantly hinders their generalization capability. To rigorously validate this, we conducted a series of empirical analyses, focusing on how a simplified, grayscale representation of images affects key aspects of deepfake detection performance. Specifically, we aimed to address two core propositions:

- **Hypothesis 1 (H1):** Removing color dependency, through grayscale processing, will *reduce the intrinsic variability between deepfake datasets* (i.e., real and synthetic images from different sources) and *enhance the separability* between real and generated images within these datasets, thereby improving *cross-dataset generalization*.
- **Hypothesis 2 (H2):** Removing color dependency will *improve the separability of real* and generated images across various individual generative models, leading to better *cross-generator generalization*.

To conduct this foundational analysis, we employed a standard grayscale conversion based on the BT.601 luminance formula as a control mechanism. For an input RGB image $\mathbf{x} = [\mathbf{x}^{(R)}, \mathbf{x}^{(G)}, \mathbf{x}^{(B)}]$, where $\mathbf{x}^{(R)}, \mathbf{x}^{(G)}$, and $\mathbf{x}^{(B)}$ represent the red, green, and blue channels respectively, the grayscale image \mathbf{x}_{gray} is computed as:

$$\mathbf{x}_{\text{gray}} = 0.299 \cdot \mathbf{x}^{(R)} + 0.587 \cdot \mathbf{x}^{(G)} + 0.114 \cdot \mathbf{x}^{(B)}$$
 (1)

This initial grayscale conversion serves to isolate the impact of color information, allowing us to observe if its removal fundamentally benefits deepfake distinction.

Verifying Hypothesis 1 (H1). To assess the impact on cross-dataset variability and separability, we analyzed the distributions of original RGB and grayscale images from three diverse datasets: DiffusionForensics Wang et al. (2023), GenImage Zhu et al. (2024), and COCO Lin et al. (2014). We utilized UMAP (Uniform Manifold Approximation and Projection) McInnes et al. (2018) to map high-dimensional image features into a two-dimensional space for visualization. As depicted in Figure 2 (a) and (b), our observations revealed that after grayscale processing, the feature distributions of these datasets became noticeably more aligned. This suggests that the elimination of color dependency reduces the inherent inter-dataset variability. More importantly, to quantify the separability between real and generated images, we employed two metrics: Density Overlap Li et al. (2021) and Maximum Mean Discrepancy (MMD) Cheng & Xie (2021). Lower density overlap and higher MMD values indicate better separation. Our results (illustrated in Figure 2 (c) and (d), corresponding to real/generated image separation) unequivocally demonstrated that grayscale images significantly increased the separation between real and generated image distributions. This finding strongly supports H1, indicating that color often introduces distracting information that impedes effective cross-dataset generalization.

Verifying Hypothesis 2 (H2). To investigate the effect on cross-generator generalization, we conducted similar experiments by analyzing image distributions from various generative models. Focusing on the DDPM Ho et al. (2020) generator as a representative example, we compared the

UMAP distributions of real and generated images before and after grayscale processing. As shown in Figure 2 (e) and (f), grayscale conversion led to a more pronounced separation between real and generated image distributions for this specific generator, consistently indicated by lower Density Overlap and higher MMD scores. This consistent improvement across different generators further substantiates H2, highlighting that generative models leave distinct, non-color-dependent artifacts that become more salient when color information is removed.

Our empirical analysis strongly suggests that excessive reliance on color dependency within deep-fake detection models significantly compromises their generalization capability. Grayscale processing, by mitigating this dependency, reveals more generalizable distinguishing features between real and generated images, laying the groundwork for improved detection across diverse scenarios.

3.2 Adaptive Gray (AG) Framework

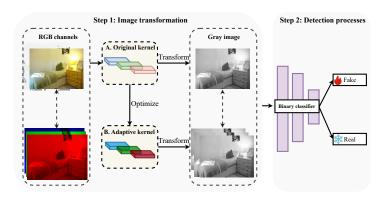


Figure 3: Pipeline of AG Training. AG is optimized from the original grayscale (OG) parameters. The adaptive training process iteratively adjusts both the AG kernel and binary classifier to improve detection performance.

Inspired by the empirical evidence that color dependency significantly impedes deepfake detection generalization, we propose *Adaptive Gray* (AG). Unlike conventional methods that rely on fixed color-to-grayscale conversion coefficients (as used in our empirical analysis), AG introduces a novel, data-driven approach to learn the optimal linear combination of RGB channels. This adaptive transformation is designed to effectively suppress misleading color information while accentuating the subtle, intrinsic textural artifacts that are highly indicative of synthetic origins.

The core idea of AG is to allow the deepfake classifier to actively participate in defining the most discriminative grayscale representation. As intuitively illustrated in Figure 1, AG processing amplifies these subtle texture-based artifacts that are typically obscured by color variations in generated images, thereby making the distinguishing features more salient for the classifier. The "Adaptive" nature of our framework stems from an adaptive training process that jointly optimizes both the image transformation (grayscale conversion parameters) and the subsequent detection stages.

3.2.1 FORMALIZING THE ADAPTIVE GRAYSCALE TRANSFORMATION

We define the adaptive grayscale function, $G(\mathbf{x}; \mathbf{w})$, which transforms an input RGB image $\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$ (with height H, width W, and 3 color channels) into a single-channel grayscale image $\mathbf{x}' \in \mathbb{R}^{H \times W}$. This transformation is a linear combination of the input image's red, green, and blue channels:

$$\mathbf{x}' = \mathbf{G}(\mathbf{x}; \mathbf{w}) = w_B \cdot \mathbf{x}^{(R)} + w_G \cdot \mathbf{x}^{(G)} + w_B \cdot \mathbf{x}^{(B)}$$
(2)

Here, $\mathbf{w} = [w_R, w_G, w_B]$ represents the learnable parameters (weights) for the grayscale transformation. Unlike the fixed coefficients in standard grayscale methods like BT.601 (Eq. 1), these parameters w_R, w_G, w_B are not predetermined. Instead, they are dynamically optimized during the training process to best suit the deepfake detection task. The resulting image \mathbf{x}' is what we refer to as the *Adaptive Gray* (AG) image.

3.2.2 CO-ADAPTIVE TRAINING ALGORITHM

The effectiveness of AG lies in its co-adaptive training process, where the grayscale transformation parameters \mathbf{w} and the deepfake binary classifier parameters θ are optimized in an alternating manner. This allows the grayscale conversion to adapt to the classifier's needs, and vice-versa, fostering a representation that maximizes separability between real and fake images. The training pipeline is illustrated in Figure 3, and the adaptive training process can be formally expressed as follows:

Let L(y, t) be the classification loss function (e.g., Binary Cross-Entropy), where y is the predicted probability from the classifier and t is the true binary label (real or fake). Let $f(\cdot; \theta)$ denote the deepfake binary classifier with parameters θ .

Step 1: Optimize Classifier Parameters (θ) In this step, the AG parameters w are held fixed. The classifier f is then trained on the images already transformed by the current w, with the objective of minimizing the classification loss:

$$\min_{\theta} L(\mathbf{f}(\mathbf{G}(\mathbf{x}; \mathbf{w}); \theta), \mathbf{t})$$
 (3)

This step enables the classifier to learn how to best classify images given the current adaptive grayscale representation, adapting its features to the information exposed by G(x; w).

Step 2: Optimize Adaptive Grayscale Parameters (w) Subsequently, the classifier parameters θ are fixed. The AG parameters w are then optimized to further minimize the same classification loss, effectively learning how to transform the input images to make them most discriminative for the fixed classifier:

$$\min_{\mathbf{w}} L(\mathbf{f}(\mathbf{G}(\mathbf{x}; \mathbf{w}); \theta), \mathbf{t}) \tag{4}$$

This crucial step allows the AG kernel to dynamically adjust its channel weights, pushing the grayscale transformation towards a representation that highlights the most salient generative artifacts, thereby improving the overall model's generalization without requiring complex architectures or large additional parameters.

This alternating optimization process, typically performed over several training cycles, ensures that both the image representation and the classification model co-evolve, leading to a fine-tuned system that is highly effective and generalizable for deepfake detection.

4 EXPERIMENT SETUP

In this section, we introduce the experiment setup, including the datasets, models, baselines, metrics and implementation details.

Train Datasets. Following Wang et al. (2023), we use DiffusionForensics for training. Real images come from the *bedroom* category of LSUN (Yu et al., 2015). In total, the dataset contains 40,000 real images in the training set. Generated images are produced by ADM (Dhariwal & Nichol, 2021) in the same category (40k).

Test Datasets. Following previous studies Wang et al. (2023); Zhu et al. (2024); Tan et al. (2024), we employed two types of test datasets in our study. (1) The first is the test dataset from *LSUN_bedroom* subset of the DiffusionForensics Wang et al. (2023) dataset Wang et al. (2023), which serves as an in-distribution dataset. This test set includes 1,000 real bedroom images, along with 1,000 generated images from each of the various generators trained on this dataset, including ADM Dhariwal & Nichol (2021), DDPM Ho et al. (2020), IDDPM Nichol & Dhariwal (2021), IF Saharia et al. (2022), Midjourney Mid (2022), PNDM Liu et al. (2022), ProGAN Karras et al. (2018), SD Rombach et al. (2021), and VQDM Gu et al. (2022). (2) The second is the GenImage dataset Zhu et al. (2024), used to evaluate cross-dataset generalization capability. GenImage Zhu et al. (2024) is specifically designed to assess the generalization performance of deepfake detection models and includes 6,000 real images and 6,000 high-quality generated images, covering 1,000 image categories for each of the generative models. Notably, it features images generated by various SOTA diffusion models, including Wukong Wuk (2022), Midjourney Mid (2022), VQDM Gu et al. (2022), and Stable Diffusion Rombach et al. (2021).

Table 1: The in-distribution performance of deepfake tasks. **Bold** and <u>underline</u> indicate the best and the second-best performance, respectively.

| Methods | In-distribution Performance | | | | | | |
|----------|-----------------------------|--------|--------|--|--|--|--|
| Michious | ACC | AP | TPR | | | | |
| ResNet | 90.05 | 93.26 | 96.70 | | | | |
| DIRE | 95.10 | 99.95 | 99.70 | | | | |
| NPR | 88.60 | 99.80 | 97.80 | | | | |
| DE-FAKE | 91.45 | 94.28 | 80.30 | | | | |
| OG | 99.95 | 99.99 | 99.99 | | | | |
| AG | 99.95 | 100.00 | 100.00 | | | | |

Baselines. In this study, we primarily compare our method against three baselines: NPR Tan et al. (2024), DIRE Wang et al. (2023), and DE-FAKE Sha et al. (2023), which represent state-of-the-art approaches in deepfake detection. These three methods are among the latest and most advanced techniques, each demonstrating high detection performance in their respective experiments. Unless otherwise noted, we use authors' official code/checkpoints or faithful re-implementations with the hyperparameters reported in the original papers to ensure fairness.

Metrics. We selected evaluation metrics that align with previous research Wang et al. (2020); Tan et al. (2024); Ricker et al. (2024); Qian et al. (2020); Sinitsa & Fried (2024) while considering real-world applicability. Common metrics such as Average Precision (AP) Wang et al. (2020); Tan et al. (2024); Ricker et al. (2024) and Average Accuracy (ACC) Tan et al. (2024); Qian et al. (2020); Sinitsa & Fried (2024) provide a comprehensive view of classifier performance. While AP evaluates performance across varying thresholds, real-world applications often lack prior knowledge of critical factors such as the true label of the image, the generative model used, and specific requirements (e.g., prioritizing either minimizing false positives or false negatives). As noted by Carlini Carlini et al. (2022) and other researchers Ho et al. (2017); Kantchelian et al. (2015); Kolter & Maloof (2006), evaluating models at lower False Positive Rates (FPR) provides a more realistic assessment. In line with Ricker Ricker et al. (2024), we included the *True Positive Rate (TPR)* at a fixed FPR, setting FPR to 5% to balance sensitivity with the minimization of false positives, ensuring a practical threshold for generative image detection. Threshold calibration and scoring. We select a single global decision threshold on a held-out validation split from the **training domain** (LSUN_bedroom, ADM) to achieve 5% FPR. This fixed threshold is then used to report TPR@5%FPR and ACC on all test sets (in-distribution, cross-generator, cross-dataset). AP is threshold-free.

Implementation details. We used ResNet50 He et al. (2016) as the backbone model for our classifier. We optimized the training process using the Adam optimizer with a learning rate of 10^{-4} . The batch size was set to 32, and training was conducted over 400 epochs, with the best model evaluated. For image preprocessing, we resized images to 224×224 pixels without additional data augmentation. Unlike prior work Wang et al. (2020), which shows that data augmentation can improve performance, we omitted it to directly assess the method's inherent generalization. All experiments were run on an Ubuntu server with an NVIDIA GeForce RTX 2080 Ti GPU.

5 Evaluation

In this section, we present comprehensive experimental results that empirically validate our central hypothesis regarding the detrimental impact of color dependency on deepfake detection generalization, and demonstrate the superior efficacy of our Adaptive Gray (AG) method in addressing this challenge. We detail our experimental setup, followed by an in-depth analysis of AG's performance across various crucial aspects: in-distribution performance, cross-generator generalization, cross-dataset generalization, robustness to unseen perturbations, inference efficiency, an ablation study clarifying the role of adaptive learning, and a qualitative analysis.

5.1 In-distribution Performance

For fair evaluation, all models were trained and tested on in-distribution data from the same generative model (ADM Dhariwal & Nichol (2021)) and image category (LSUN_bedroom Yu et al.

Table 2: The cross-generator performance of deepfake tasks. Results demonstrate the superior generalization capabilities of AG and OG methods.

| Methods | DDPM | | IDDPM | | IF | | Midjourney | | Methods | Mean | | | | | | |
|---------|-------|-------|-------|--------|--------------|--------------|------------|-------|---------|-------|--------------|-------|---------|-------------|--------------|--------------|
| Methods | ACC | AP | TPR | ACC | AP | TPR | ACC | AP | TPR | ACC | AP | TPR | Methods | ACC | AP | TPR |
| ResNet | 56.56 | 48.81 | 11.07 | 50.00 | 51.11 | 5.7 | 50.00 | 38.74 | 0 | 92.91 | 94.40 | 98.00 | ResNet | 56.18 | 8 55.18 | 17.95 |
| DIRE | 93.72 | 99.26 | 98.60 | 94.90 | 99.74 | 99.20 | 50.00 | 47.85 | 5.10 | 91.18 | 90.91 | 89.50 | | 30.16 | | |
| NPR | 93.97 | 99.93 | 98.80 | 91.80 | 99.80 | 99.40 | 92.67 | 98.10 | 99.75 | 91.60 | 91.90 | 97.20 | DIRE | 82.67 87.66 | 97.66 | 74.68 |
| DE-FAKE | 90.16 | 92.31 | 77.73 | 93.20 | 96.33 | 99.96 | 72.20 | 82.08 | 28.50 | 83.27 | 13.62 | 2.00 | | | 67.00 | 74.00 |
| OG | 99.32 | 99.96 | 99.74 | 99.45 | 99.99 | 99.80 | 66.55 | 88.24 | 54.30 | 90.09 | 6.06 | 0 | NPR | 91.21 | <u>97.92</u> | <u>95.19</u> |
| AG | 99.96 | 100 | 99.94 | 99.95 | 100 | 99.90 | 97.40 | 99.99 | 99.90 | 94.91 | 98.02 | 99.00 | | | | |
| Methods | PNDM | | | ProGAN | | | SD | | VQDM | | DE-FAKE | 76.96 | 74.68 | 42.46 | | |
| | ACC | AP | TPR | ACC | AP | TPR | ACC | AP | TPR | ACC | AP | TPR | DL-TAKE | 70.90 | 74.00 | 72.70 |
| ResNet | 50.00 | 42.55 | 3.50 | 50.00 | 54.93 | 10.50 | 50.00 | 69.31 | 14.40 | 50.00 | 41.65 | 0.40 | | | | |
| DIRE | 92.45 | 97.87 | 89.70 | 94.00 | <u>98.78</u> | <u>95.90</u> | 50.00 | 66.99 | 20.20 | 95.10 | <u>99.87</u> | 99.30 | OG | 85.28 | 83.42 | 18.07 |
| NPR | 91.20 | 99.99 | 99.90 | 76.67 | 95.65 | 67.50 | 93.80 | 98.40 | 99.30 | 98.00 | 99.57 | 99.70 | | | | |
| DE-FAKE | 88.25 | 93.14 | 73.10 | 78.15 | 86.14 | 41.70 | 55.15 | 67.96 | 7.70 | 55.30 | 65.89 | 9.00 | | | | |
| OG | 99.00 | 99.96 | 99.80 | 85.70 | 98.28 | 91.00 | 57.15 | 76.98 | 29.40 | 84.95 | 97.89 | 87.90 | AG | 97.10 | 99.73 | 99.71 |
| AG | 99.90 | 100 | 99.70 | 97.85 | 99.98 | 99.90 | 87.35 | 99.88 | 99.90 | 99.45 | 99.95 | 99.40 | | | | |

Table 3: The cross-dataset performance of deepfake tasks. AG demonstrates superior generalization across diverse image categories and generative models.

| Methods | Wukong | | Midjourney | | VQDM | | SD | | Mean | | | | | | |
|---------|--------|--------------|------------|-------|-------|-------|--------------|--------------|-------|-------|--------------|-------|-------|-------|-------|
| Methods | ACC | AP | TPR | ACC | AP | TPR | ACC | AP | TPR | ACC | AP | TPR | ACC | AP | TPR |
| ResNet | 41.08 | 41.69 | 0.12 | 48.15 | 44.95 | 1.65 | 51.22 | 51.20 | 6.45 | 46.64 | 40.60 | 0.08 | 46.77 | 44.61 | 2.08 |
| DIRE | 53.57 | 70.12 | 12.85 | 53.43 | 70.56 | 11.97 | 54.10 | 68.92 | 15.65 | 54.75 | 76.05 | 23.46 | 53.96 | 71.41 | 15.98 |
| NPR | 50.00 | 55.11 | 11.33 | 50.00 | 46.75 | 2.93 | 50.00 | 45.87 | 4.97 | 50.00 | 59.11 | 11.49 | 50.00 | 51.71 | 7.68 |
| DE-FAKE | 50.20 | 52.19 | 5.48 | 52.63 | 57.61 | 8.30 | 61.18 | 68.99 | 23.35 | 51.16 | 53.53 | 6.33 | 53.79 | 58.08 | 10.87 |
| OG | 49.78 | 51.29 | 4.87 | 54.83 | 60.85 | 14.15 | 63.40 | 76.13 | 25.15 | 51.82 | 57.37 | 8.45 | 54.96 | 61.41 | 13.16 |
| AG | 54.92 | <u>66.20</u> | 16.85 | 64.33 | 78.60 | 36.65 | <u>59.78</u> | <u>75.99</u> | 28.47 | 60.57 | <u>74.19</u> | 29.15 | 59.90 | 73.75 | 27.78 |

(2015)). As shown in Table 1, our AG method achieves outstanding in-distribution performance (ACC 99.95%, AP 100.00%, TPR 100.00%), closely followed by the Original Gray (OG) method (ACC 99.95%, AP 99.99%, TPR 99.99%). While most compared methods show high accuracy in this setting, our results establish a strong baseline for AG's foundational capabilities before assessing its critical generalization performance.

5.2 Cross-generator Performance

We conducted cross-generator experiments to analyze generalization capabilities. Models, trained on ADM-generated and LSUN_bedroom real images, were evaluated on various **unseen** generative models within the DiffusionForensics dataset Wang et al. (2023) (fixed 'bedroom' category).

Table 2 strikingly validates our Hypothesis 2 (H2) from Section 3.1. Our OG method, applying simple fixed grayscale transformation, demonstrates a substantial leap in cross-generator generalization over baseline ResNet50 and even some SOTA methods. This supports our assertion that mitigating color dependency significantly enhances generalization to novel generators. Building on this, our AG method further elevates performance, consistently achieving the highest average ACC, AP, and TPR values. AG is uniquely the only method where its mean ACC, AP, and TPR all exceed 97%, often reaching 100% accuracy for individual generators. This superior performance of AG over OG underscores the effectiveness of adaptively learning optimal grayscale coefficients to best isolate generative artifacts.

In contrast, SOTA methods like DIRE Wang et al. (2023), DE-FAKE Sha et al. (2023), and NPR Tan et al. (2024) show varied cross-generator performance. Their fluctuations, contrasted with the consistent gains from OG and AG, highlight the persistent challenge of generalization when relying heavily on color-rich feature spaces. Our findings strongly suggest that addressing color dependency is a crucial step towards robust cross-generator detection.

5.3 Cross-dataset Performance

Building on cross-generator insights, we conducted cross-dataset testing to rigorously evaluate generalization under more challenging conditions. This involved diversifying test sets across both image categories (e.g., 'bedroom' to faces, landscapes) and generative models, exposing classifiers to a broader range of real-world data variations.

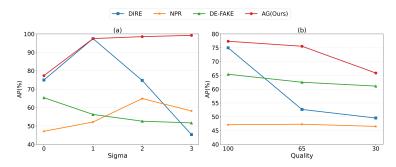


Figure 4: Robustness to unseen perturbations. Figure (a) illustrates performance against Gaussian blur, and Figure (b) shows performance against JPEG compression for AG, NPR, DE-FAKE, and DIRE. We evaluate the average AP of all methods on GenImage, using images from BigGAN, Midjourney, Stable Diffusion, and Wukong models.

Table 3 emphatically confirms Hypothesis 1 (H1) from Section 3.1 and highlights AG's superior generalization. The OG method, by simply removing fixed color dependency, again shows notable improvement over baselines and SOTA methods, corroborating our finding that color information significantly impedes cross-dataset generalization. Crucially, our AG method consistently achieves the highest overall results in this demanding scenario, with average ACC of **59.90**%, AP of **73.75**%, and TPR of **27.78**%. These represent substantial relative improvements (up to 19.9% ACC, 22.04% AP, and 20.10% TPR) compared to the best SOTA methods, underscoring AG's remarkable robustness to diverse content and models.

In summary, AG consistently outperforms or remains highly competitive with other SOTA detection methods Wang et al. (2023); Tan et al. (2024); Sha et al. (2023) across all testing scenarios—indistribution, cross-generator, and especially the challenging cross-dataset conditions. These results not only demonstrate AG's superior practical utility but also provide strong empirical validation for our core hypothesis that mitigating color dependency is fundamental for achieving robust and generalizable deepfake detection.

5.4 ROBUSTNESS TO UNSEEN PERTURBATIONS

Robustness to common image perturbations is crucial for real-world deepfake detection, as images often undergo degradations. We evaluated AG against DIRE Wang et al. (2023), NPR Tan et al. (2024), and DE-FAKE Sha et al. (2023) under Gaussian blur ($\sigma=1,2,3$) and JPEG compression (quality 65, 30), following previous studies Wang et al. (2023).

As shown in Figure 4, AG consistently outperforms baselines across all blur and compression levels. For Gaussian blur, AG maintains high and stable Average Precision (AP). Under JPEG compression, AG exhibits significantly less performance degradation; specifically, while DIRE experiences a substantial AP drop at JPEG quality 30, AG remains remarkably resilient. This outcome reinforces our core hypothesis: AG's adaptive grayscale processing reduces reliance on superficial color information, making it inherently more resilient to common image degradations that often corrupt color channels. Its focus on robust, underlying texture-based artifacts is key to this enhanced robustness.

6 Conclusion

In this work, we tackled the challenge of enhancing binary deepfake classifiers' generalization by reducing color dependency in real vs. generated image detection. We proposed that color discrepancies can hinder detection accuracy, leading to the development of our grayscale processing framework, **Adaptive Gray (AG)**. Through adaptive training of both grayscale parameters and the classifier, AG demonstrated superior generalization across datasets and testing conditions, outperforming SOTA methods. Our findings suggest that focusing on texture over color can improve detection resilience, offering a promising direction for generalizable detection systems applications.

REFERENCES

486

487

497

498 499

500

501

502

504

505

507

508

509

510 511

512

513

514 515

516

517

518

519

520 521

522

523 524

525

526

527

528

529

530 531

532

- 488 Midjourney. https://www.midjourney.com/home, 2022.
- Wukong. https://xihe.mindspore.cn/modelzoo/wukong, 2022.
- strengthens Minns labor government protections against deepfakes and 491 image-based abuse. **NSW** Government ministerial release. 09 492 2025. **URL** https://www.nsw.gov.au/ministerial-releases/ 493 minns-labor-government-strengthens-protections-against-deepfakes-and-image-based-a
- Amendments to the Crimes Act 1900 criminalising sexually explicit deepfakes; penalties up to three years.
 - Andrew Brock. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint* arXiv:1809.11096, 2018.
 - Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. Membership inference attacks from first principles. In 2022 IEEE Symposium on Security and Privacy (SP), pp. 1897–1914. IEEE, 2022.
 - Liang Chen, Yong Zhang, Yibing Song, Lingqiao Liu, and Jue Wang. Self-supervised learning of adversarial example: Towards good generalizations for deepfake detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 18710–18719, 2022.
 - Xiuyuan Cheng and Yao Xie. Neural tangent kernel maximum mean discrepancy. *Advances in Neural Information Processing Systems*, 34:6658–6670, 2021.
 - Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
 - Joel Frank, Franziska Herbert, Jonas Ricker, Lea Schönherr, Thorsten Eisenhofer, Asja Fischer, Markus Dürmuth, and Thorsten Holz. A representative study on human detection of artificially generated media across countries. In 2024 IEEE Symposium on Security and Privacy (SP), pp. 55–73. IEEE, 2024.
 - Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10696–10706, 2022.
 - Yuanfang Guo, Xiaochun Cao, Wei Zhang, and Rui Wang. Fake colorized image detection. *IEEE Transactions on Information Forensics and Security*, 13(8):1932–1944, 2018.
 - Alexandros Haliassos, Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Lips don't lie: A generalisable and robust approach to face forgery detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5039–5049, 2021.
 - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
 - Grant Ho, Aashish Sharma, Mobin Javed, Vern Paxson, and David Wagner. Detecting credential spearphishing in enterprise settings. In *26th USENIX security symposium (USENIX security 17)*, pp. 469–485, 2017.
 - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Yonghyun Jeong, Doyeon Kim, Youngmin Ro, Pyounggeon Kim, and Jongwon Choi. Fingerprintnet: Synthesized fingerprints for generated image detection. In *European Conference on Computer Vision*, pp. 76–94. Springer, 2022.
- Alex Kantchelian, Michael Carl Tschantz, Sadia Afroz, Brad Miller, Vaishaal Shankar, Rekha Bachwani, Anthony D Joseph, and J Doug Tygar. Better malware ground truth: Techniques for weighting anti-virus vendor labels. In *Proceedings of the 8th ACM Workshop on Artificial Intelligence and Security*, pp. 45–56, 2015.

- Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *ICLR Workshop*, 2018.
 - Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4401–4410, 2019.
 - J Zico Kolter and Marcus A Maloof. Learning to detect and classify malicious executables in the wild. *Journal of Machine Learning Research*, 7(12), 2006.
 - Zhuang Li, Jingyan Qin, Xiaotong Zhang, and Yadong Wan. Addressing class overlap under imbalanced distribution: An improved method and two metrics. *Symmetry*, 13(9):1649, 2021.
 - Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
 - Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778*, 2022.
 - Francesco Marra, Diego Gragnaniello, Luisa Verdoliva, and Giovanni Poggi. Do gans leave artificial fingerprints? In 2019 IEEE conference on multimedia information processing and retrieval (MIPR), pp. 506–511. IEEE, 2019.
 - Scott McCloskey and Michael Albright. Detecting gan-generated imagery using color cues. *arXiv* preprint arXiv:1812.08247, 2018.
 - Scott McCloskey and Michael Albright. Detecting gan-generated imagery using saturation cues. In 2019 IEEE international conference on image processing (ICIP), pp. 4584–4588. IEEE, 2019.
 - Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
 - Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.
 - Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in frequency: Face forgery detection by mining frequency-aware clues. In *European conference on computer vision*, pp. 86–103. Springer, 2020.
 - Jonas Ricker, Denis Lukovnikov, and Asja Fischer. Aeroblade: Training-free detection of latent diffusion images using autoencoder reconstruction error. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9130–9140, June 2024.
 - Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.
 - Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1–11, 2019.
 - Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
 - Ayush Sarkar, Hanlin Mai, Amitabh Mahapatra, Svetlana Lazebnik, David A Forsyth, and Anand Bhattad. Shadows don't lie and lines can't bend! generative models don't know projective geometry... for now. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 28140–28149, 2024.
 - Axel Sauer, Kashyap Chitta, Jens Müller, and Andreas Geiger. Projected gans converge faster. Advances in Neural Information Processing Systems, 34:17480–17492, 2021.

Zeyang Sha, Zheng Li, Ning Yu, and Yang Zhang. De-fake: Detection and attribution of fake images generated by text-to-image generation models. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pp. 3418–3432, 2023.

Sergey Sinitsa and Ohad Fried. Deep image fingerprint: Towards low budget synthetic image detection and model lineage analysis. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 4067–4076, 2024.

Chuangchuang Tan, Yao Zhao, Shikui Wei, Guanghua Gu, Ping Liu, and Yunchao Wei. Rethinking the up-sampling operations in cnn-based generative network for generalizable deepfake detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 28130–28139, 2024.

Chengrui Wang and Weihong Deng. Representative forgery mining for fake face detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 14923–14932, 2021.

Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8695–8704, 2020.

Zhendong Wang, Jianmin Bao, Wengang Zhou, Weilun Wang, Hezhen Hu, Hong Chen, and Houqiang Li. Dire for diffusion-generated image detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22445–22455, 2023.

Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv* preprint arXiv:1506.03365, 2015.

Ning Yu, Larry S Davis, and Mario Fritz. Attributing fake images to gans: Learning and analyzing gan fingerprints. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 7556–7566, 2019.

Mingjian Zhu, Hanting Chen, Qiangyu Yan, Xudong Huang, Guanyu Lin, Wei Li, Zhijun Tu, Hailin Hu, Jie Hu, and Yunhe Wang. Genimage: A million-scale benchmark for detecting ai-generated image. *Advances in Neural Information Processing Systems*, 36, 2024.

A APPENDIX

A.1 EFFICIENCY EVALUATION

Table 4: Inference Time Comparison of Different Methods.

| Method | Average Inference Time per Image (s) |
|-----------|--------------------------------------|
| DIRE | 269.50 |
| DE-FAKE | 12.12 |
| NPR | 8.66×10^{-3} |
| AG (Ours) | 6.12×10^{-3} |

Our AG method significantly outperforms most contemporary approaches in inference speed. As Table 4 shows, AG achieves the fastest average inference time at 6.12×10^{-3} seconds per image, closely followed by NPR (8.66×10^{-3} s). In contrast, DE-FAKE takes 12.12 seconds, and DIRE incurs a notably long 269.50 seconds per image.

AG's remarkable efficiency stems from its simple, linear processing during inference. Once w is learned, image transformation involves only straightforward multiplication and addition, approximating the speed of a standard ResNet50 classifier He et al. (2016). This contrasts sharply with DE-FAKE's multi-component pipeline (CLIP encoders, large classifier, optional BLIP for captioning) and DIRE's computationally intensive diffusion model reconstruction, both of which are impractical for real-time or large-scale applications.

Table 5: Cross-generator performance comparison including AG + DIRE. **Bold** and <u>underline</u> denote the best and second-best utility, respectively.

| Methods | Cross-generator Performance | | | | | | |
|-----------|-----------------------------|--------------|-------|--|--|--|--|
| Methods | ACC | AP | TPR | | | | |
| ResNet | 56.18 | 55.18 | 17.95 | | | | |
| OG | 85.28 | 83.42 | 18.07 | | | | |
| DIRE + AG | 92.76 | <u>97.51</u> | 90.47 | | | | |
| AG | 97.10 | 99.73 | 99.71 | | | | |

Table 6: Cross-dataset performance comparison including AG + DIRE.

| Methods | Cross-dataset Performance | | | | | | |
|-----------|---------------------------|-------|-------|--|--|--|--|
| Michiods | ACC | AP | TPR | | | | |
| ResNet | 46.77 | 44.61 | 2.08 | | | | |
| OG | 54.96 | 61.41 | 13.16 | | | | |
| DIRE + AG | 58.74 | 67.23 | 0 | | | | |
| AG | 59.90 | 73.75 | 27.78 | | | | |

Combined with its strong generalization, AG's exceptional efficiency positions it as a practical, deployable solution for real-time deepfake detection and high-volume image analysis, alleviating common computational bottlenecks.

A.2 ABLATION STUDY: THE IMPACT OF ADAPTIVE GRAYSCALE LEARNING

Our methodology posits that both basic grayscale processing and adaptive learning of its coefficients enhance deepfake detection generalization. To dissect these contributions, we address two research questions (RQs) related to our empirical hypotheses (H1 and H2) from Section 3.1:

- **RQ1:** Basic Grayscale Efficacy Does fixed grayscale compression improve generalization by retaining relevant features?
- **RQ2:** Value of Adaptive Learning: Does optimizing grayscale parameters (AG) further enhance generalization beyond fixed grayscale, suggesting a more discriminative compression learned by machines?

To verify **RQ1**, we evaluated the Original Gray (OG) method, which uses standard BT.601 grayscale conversion (Eq. 1) without adaptive training. As shown in Table 2 and Table 3, even fixed OG processing notably improved generalization. Specifically, in cross-generator testing, OG achieved approximate increases of 29% in ACC and 28% in AP over baseline ResNet50 He et al. (2016). In cross-dataset testing, OG still gained about 8% in ACC and 16% in AP compared to the baseline. These findings strongly support RQ1, confirming that simply mitigating color dependency through grayscale inherently enhances generalization by retaining critical texture-based features.

To verify **RQ2**, we compared OG and AG to see if optimizing grayscale parameters further enhances generalization. This addresses whether machine learning can discover a more optimal grayscale kernel. Results from Table 2 and Table 3 unequivocally confirm that adaptively training AG parameters substantially improves classifier generalization beyond fixed OG. In cross-generator testing, AG surpassed OG with approximate increases of 12% in ACC, 16% in AP, and a remarkable 82% in TPR. Similarly, for cross-dataset testing, AG demonstrated superior generalization with gains of about 5% in ACC, 12% in AP, and 15% in TPR over OG. These results strongly support RQ2, indicating that the co-adaptive training process enables AG to learn a more discriminative representation by optimizing the grayscale kernel for more accurate deepfake distinction.

A.3 EXPLORATORY ANALYSIS: INTEGRATING ADAPTIVE GRAY WITH DIRE

We explored combining Adaptive Gray (AG) with DIRE Wang et al. (2023), a method leveraging diffusion model reconstruction error to distinguish real from generated images. Real images typically show larger reconstruction errors.

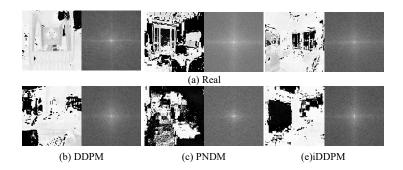


Figure 5: Frequency-domain analysis using Fast Fourier Transform (FFT) on AG grayscale images, showcasing distinct characteristics between real and generated images. (a) Real images typically demonstrate a richer frequency spectrum with organized noise patterns aligned with natural object structures. In contrast, generated images from models such as DDPM and PNDM exhibit more irregular and artificial frequency distributions. This comparison highlights that FFT-based analysis on AG-transformed images effectively enhances the separability of real from generated content in deepfake detection, by making subtle generative artifacts more visually apparent in the frequency domain.

However, as Tables 5 and 6 show, integrating AG with DIRE did not enhance generalization; instead, it led to increased instability and performance degradation. In cross-generator tests, DIRE+AG decreased average ACC, AP, and TPR by approximately 6%, 2%, and 9% respectively, compared to AG alone. Performance dropped even more sharply in cross-dataset tests, with TPR plummeting to zero at FPR=5%, indicating complete failure in that demanding scenario.

We hypothesize this arises from a fundamental incompatibility in information processing. DIRE's reconstruction already significantly reduces image content, losing shape, texture, and color. Applying AG's grayscale compression to these already feature-reduced residuals likely compounds this effect, excessively "over-filtering" meaningful features and diminishing overall performance. This suggests that while both methods effectively identify generative artifacts individually, their sequential application can be detrimental.

A.4 QUALITATIVE ANALYSIS OF AG

To further illuminate AG's underlying mechanisms and properties, we performed a qualitative frequency-domain analysis using Fast Fourier Transform (FFT) on AG-transformed images, following our quantitative results confirming AG's effectiveness through reduced color dependency. This analysis aimed to visually and analytically explain how AG enhances discriminative power.

As illustrated in Figure 5, AG processing reveals distinct low-level frequency characteristics. Real images, after AG transformation, exhibit a richer, more organized frequency spectrum with noise patterns aligning with natural structures. Conversely, generated images (e.g., ADM, DDPM) display irregular, chaotic, or unnatural frequency distributions. These irregularities are subtle generative artifacts that become significantly more pronounced and detectable in the grayscale frequency domain, often obscured by complex color patterns in RGB space. This strongly supports our hypothesis: AG effectively accentuates intrinsic, non-color-dependent artifacts for generalized deepfake detection.