# EMPIRICAL ROBUSTNESS OF PIXEL DIFFUSION UNDER-MINES ADVERSARIAL PERTURBATION AS PROTECTION AGAINST DIFFUSION-BASED MIMICRY

**Anonymous authors**Paper under double-blind review

#### **ABSTRACT**

Diffusion models have demonstrated impressive abilities in image editing and imitation, raising growing concerns about the protection of private property. A common defense strategy is to apply adversarial perturbations that can mislead a diffusion model into generating bad-quality images. However, existing research has almost entirely focused on latent diffusion models while overlooking pixel-space diffusion models. Through extensive experiments, we show that nearly all attacks designed for latent diffusion models, as well as adaptive attacks aimed at pixelspace diffusion models, fail to compromise the latter. Our analysis suggests that the weakness of latent diffusion models arises mainly from their encoder, whereas pixel-space diffusion models exhibit strong empirical robustness to adversarial perturbations. We further demonstrate that pixel-space diffusion models can serve as an effective purifier by removing adversarial patterns generated for latent diffusion models and preserving image integrity, which in turn allows them to bypass most existing protection schemes. These findings challenge the assumption that adversarial perturbations provide reliable protection for diffusion models and call for a reevaluation of their role as a protection mechanism.

# 1 Introduction

Generative diffusion models (DMs) (Ho et al., 2020; Song et al., 2020; Rombach et al., 2022) have achieved great success in generating images with high fidelity. However, this remarkable generative capability of diffusion models is accompanied by safety concerns (Zhang et al., 2023), especially on the unauthorized editing or imitation of personal images such as portraits or individual artworks (Andersen, 2023; Setty, 2023). Recent works (Liang et al., 2023; Shan et al., 2023; Salman et al., 2023; Xue et al., 2023; Zheng et al., 2023; Chen et al., 2024; Ahn et al., 2024; Liu et al., 2023) show that adversarial samples (adv-samples) for diffusion models can be applied as a **protection** against malicious editing. Small perturbations generated by conventional methods in adversarial machine learning (Madry et al., 2018; Goodfellow et al., 2014) can effectively fool popular diffusion models such as Stable Diffusion (Rombach et al., 2022) to produce chaotic results when an imitation attempt is made. However, a significantly overlooked aspect is that all the existing works focus on latent diffusion models (LDMs) and the pixel-space diffusion models (PDMs) are not studied, though there exists pixel diffusion models with strong image editing ability e.g. (Shonenkov et al.; Balaji et al., 2022). The key distinction between attacking LDMs and PDMs is that adversarial noise is applied directly to the latter, whereas in the former it is first scaled by the image encoder.

By investigating adv-samples for PDMs, we find the embarrassing fact that PDMs are **empirically** robust to existing attacks for diffusion model used in previous works. We conduct experiments on various LDMs or PDMs with different network architectures (e.g. U-Net (Ho et al., 2020), Transformer (Peebles and Xie, 2023)), different training datasets, and different input resolutions (e.g. 64, 256, 512). Through extensive experiments, we demonstrate that all the existing methods we tested (Liang and Wu, 2023; Zheng et al., 2023; Shan et al., 2023; Xue et al., 2023; Chen et al., 2024; Salman et al., 2023; Liang et al., 2023), developed to attack LDMs, fail to generate effective adv-samples for PDMs. Moreover, we conduct **adaptive attacks** for PDMs, applying strategies like gradient averaging, attacking the intermediate features and black-box attacks, but none of the attacks can effectively change the reverse diffusion process the way they do to fool LDMs. Some theoretical

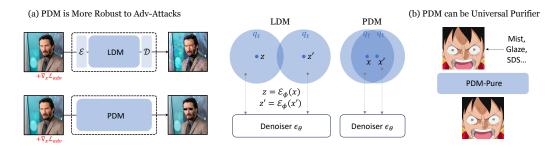


Figure 1: Overview: (a) Recent protection approaches based on adversarial perturbation against latent diffusion models (LDMs) cannot be used in pixel-space diffusion models (PDMs); The underlying reason is that the encoder of the Latent Diffusion Model (LDM) amplifies the perturbations, causing the inputs to the denoiser to have significantly different distributions. In contrast, the inputs of the PDM maintain large overlap, showing robustness. (b) A strong PDM can be used as a universal purifier to effectively remove the protective perturbation generated by existing protection methods. (Best viewed with zoom-in on a computer)



Figure 2: **PDMs Cannot be Attacked as LDMs**: LDMs can be easily fooled by running PGD to fool the denoising loss, but PDMs cannot be easily fooled when attacking the diffusion loss. DiT (Peebles and Xie, 2023) and SD (Rombach et al., 2022) are LDMs, GD (Dhariwal and Nichol, 2021) AND IF-Stage-II (Shonenkov et al.) are PDMs (Best viewed with zoom-in)

evidence for the certified robustness of the pixel-space diffusion process has been studied by (Chen et al.). Building on their results, we analyze the denoising step in the image-editing process to further support our conclusion.

Building on the insight that PDMs are strongly robust against adversarial perturbations, we further propose PDM-Pure, a universal purifier that can effectively remove the protective perturbations of different scales (e.g. Mist-v2 (Zheng et al., 2023) and Glaze (Shan et al., 2023)) based on PDMs trained in large-scale. Through extensive experiments, we demonstrate that PDM-Pure achieves way better performance than all baseline methods.

To summarize, the pixel is a barrier to adversarial attack on diffusion model (Figure 1); the diffusion process in the pixel space makes PDMs much more robust than LDMs. This property of PDMs also undermines current protection against diffusion-based mimicry because: (1) no existing attacks have proven effective in attacking PDMs, which means no protection can be achieved by fooling a PDM-based image editor, (2) all the effective protections against LDMs can be easily purified using a strong PDM.

#### 2 Preliminaries

**Generative Diffusion Models** The generative diffusion model (Ho et al., 2020; Song et al., 2020) is one type of generative model, and it has demonstrated remarkable generative capabilities in numerous

fields such as images (Rombach et al., 2022; Balaji et al., 2022), 3D data (Poole et al., 2023; Lin et al., 2022), video (Ho et al., 2022; Singer et al., 2022), stories (Pan et al., 2022; Rahman et al., 2023) and music (Mittal et al., 2021; Huang et al., 2023) generation. Diffusion models, like other generative models, are parametrized models  $p_{\theta}(\hat{x}_0)$  that can estimate an unknown distribution  $q(x_0)$ . For image generation tasks,  $q(x_0)$  is the distribution of real images.

There are two processes involved in a diffusion model, a forward diffusion process and a reverse denoising process. The forward diffusion process progressively injects noise into the clean image, and the t-th step diffusion is formulated as  $q(x_t \mid x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t \mathbf{I})$ . Accumulating the noise, we have  $q_t(x_t \mid x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_{t-1}, (1 - \bar{\alpha}_t) \mathbf{I})$ . Here  $\beta_t$  growing from 0 to 1 are pre-defined values,  $\alpha_t = 1 - \beta_t$ , and  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ . Finally,  $x_T$  will become approximately an isotropic Gaussian random variable when  $\bar{\alpha}_t \to 0$ .

Inversely,  $p_{\theta}(\hat{x}_{t-1}|\hat{x}_t)$  can generate samples from Gaussian  $\hat{x}_T \sim \mathcal{N}(0,\mathbf{I})$ , where  $p_{\theta}$  is reparameterized by learning a noise estimator  $\epsilon_{\theta}$ , the training loss is  $\mathbb{E}_{t,x_0,\epsilon}[\lambda(t)\|\epsilon_{\theta}(x_t,t)-\epsilon\|^2]$  weighted by  $\lambda(t)$ , where  $\epsilon$  is the noise used to diffuse  $x_0$  following  $q_t(x_t|x_0)$ . Finally, by iteratively applying  $p_{\theta}(\hat{x}_{t-1}|\hat{x}_t)$ , we can sample realistic images following  $p_{\theta}(\hat{x}_0)$ .

Since the above diffusion process operates directly in the pixel space, we call such diffusion models **Pixel-Space Diffusion Models (PDMs)**. Another popular choice is to move the diffusion process into the latent space to make it more scalable, resulting in the **Latent Diffusion Models (LDMs) (Rombach et al., 2022)**. More specifically, LDMs first use an encoder  $\mathcal{E}_{\phi}$  parameterized by  $\phi$  to encode  $x_0$  into a latent variable  $z_0 = \mathcal{E}_{\phi}(x_0)$ . The denoising diffusion process is the same as PDMs. At the end of the denoising process,  $\hat{z}_0$  can be projected back to the pixel space using a decoder  $\mathcal{D}_{\psi}$  parameterized by  $\psi$  as  $\hat{x}_0 = \mathcal{D}_{\psi}(\hat{z}_0)$ .

Adversarial samples (Goodfellow et al., 2014; Carlini and Wagner, 2017; Shan et al., 2023) are clean data samples perturbed by an imperceptible small noise that can fool deep neural networks into making wrong decisions. Under white-box conditions, gradient-based methods are widely used to generate adv-samples. Among them, the projected gradient descent (PGD) algorithm (Madry et al., 2018) is one of the most effective methods.

Adversarial Examples for Diffusion Models as Protection Recent works (Salman et al., 2023; Liang et al., 2023) find that adding small perturbations to clean images will make the diffusion (Ho et al., 2020; Song et al., 2020; Peebles and Xie, 2023) perform badly in noise prediction, and further generate chaotic results (Figure 1) in tasks like image editing and customized generation, which serves as one kind of **protection** against misuse of personal images. The adversarial perturbations for LDMs can be generated by optimizing the Monte-Carlo-based adversarial loss:

$$\mathcal{L}_{adv}^{\text{LDM}}(x) = \mathbb{E}_{t,\epsilon} \mathbb{E}_{z_t \sim q_t(\mathcal{E}_{\phi}(x))} \| \epsilon_{\theta}(z_t, t) - \epsilon \|_2^2.$$
 (1)

Other encoder-based losses (Shan et al., 2023; Liang and Wu, 2023; Zheng et al., 2023; Xue et al., 2023) further enhance the attack to make it more effective. With the carefully designed adversarial loss, one can run Projected Gradient Descent (PGD) (Madry et al., 2018) with  $\ell_{\infty}$  budget  $\delta$  to generate adversarial perturbations:

$$x^{k+1} = \mathcal{P}_{B_{\infty}(x^0,\delta)} \left[ x^k + \eta \operatorname{sign} \nabla_{x^k} \mathcal{L}_{adv}^{\operatorname{LDM}}(x^k) \right]$$
 (2)

In the above equation,  $\mathcal{P}_{B_\infty(x^0,\delta)}(\cdot)$  is the projection operator on the  $\ell_\infty$  ball, where  $x^0$  is the clean image to be perturbed. We use superscript  $x^k$  to represent the iterations of the PGD and subscript  $x_t$  for the diffusion steps. There are many follow-up works (Liu et al., 2023; Choi et al., 2024; Ozden et al., 2024; Kang et al., 2025; Li et al., 2024; Shen et al., 2025) trying to improve the attacks, but all of these works study only latent diffusion models.

# 3 RETHINKING ADVERSARIAL EXAMPLES FOR DIFFUSION MODELS

#### 3.1 DIFFUSION MODELS DEMONSTRATE STRONG ADVERSARIAL ROBUSTNESS

While there are many approaches that adopt adversarial perturbation to fool diffusion models, most of them focus only on latent diffusion models due to the wide impact of Stable Diffusion (Podell et al., 2023); no attempts have been made to attack PDMs. This lack of investigation may mislead us to conclude that diffusion models, like most deep neural networks, are vulnerable to adversarial perturbations. Similar to Eq 1, the diffusion loss based attack for PDM can be formulated as in the following equation, we also use PGD to update the perturbation as in Eq 2:

$$\mathcal{L}_{adv}^{\text{PDM}}(x) = \mathbb{E}_{t,\epsilon} \mathbb{E}_{x_t \sim q_t(x)} \| \epsilon_{\theta}(x_t, t) - \epsilon \|_2^2.$$
(3)

We find that PDMs are in fact robust against this form of attack (Figure 2), which means all the existing attacks against diffusion models are, in fact, special cases of attacks against the LDMs only. We conduct extensive experiments on popular LDMs and PDMs structures including Diffusion Transformer (DiT) (Peebles and Xie, 2023), Guided Diffusion (GD) (Dhariwal and Nichol, 2021), Stable Diffusion (SD) (Podell et al., 2023), and DeepFloyd (IF) (Shonenkov et al.), and demonstrate in Table 1 that only the LDMs can be attacked and PDMs are not as susceptible to adversarial perturbations: for PDMs, the image quality does not significantly decrease due to the perturbation both visually and quantitatively. More details and analysis can be found in the experiment section.

Prior to this study, there may have been a prevailing belief that diffusion models could be easily deceived. However, our research reveals an important distinction: it is the LDMs that exhibit vulnerability, while the PDMs demonstrate significantly higher adversarial robustness.

Models	FI	D-scor	e†	5	SSIM ↓			LPIPS	<b>↑</b>	IA	A-Score	÷	Type
$\delta = 4/255$	Clean	Adv	Δ	Clean	Adv	Δ	Clean	Adv	Δ	Clean	Adv	Δ	
DiT-256	131	167	+36	0.37	0.35	-0.02	0.44	0.54	+0.10	0.74	0.70	-0.04	LDM
SD-V-1.4	44	114	+70	0.68	0.55	-0.13	0.22	0.46	+0.24	0.92	0.84	-0.08	LDM
SD-V-1.5	45	113	+68	0.73	0.59	-0.14	0.20	0.38	+0.138	0.94	0.89	-0.05	LDM
GD-ImageNet	109	109	+0	0.66	0.66	-0.00	0.21	0.21	+0.00	0.90	0.90	-0.00	PDM
IF-I	186	187	+1	0.59	0.58	-0.01	0.14	0.14	+0.00	0.86	0.86	-0.00	PDM
IF-II	85	87	+2	0.84	0.84	-0.00	0.15	0.15	+0.00	0.91	0.91	-0.00	PDM
$\delta = 8/255$	Clean	Adv	Δ	Clean	Adv	Δ	Clean	Adv	Δ	Clean	Adv	Δ	
DiT-256	131	186	+55	0.37	0.31	-0.06	0.44	0.63	+0.19	0.74	0.66	-0.08	LDM
SD-V-1.4	44	178	+134	0.68	0.44	-0.24	0.22	0.60	+0.38	0.92	0.78	-0.14	LDM
SD-V-1.5	45	179	+134	0.73	0.49	-0.24	0.20	0.51	+0.31	0.94	0.84	-0.10	LDM
GD-ImageNet	109	110	+1	0.66	0.64	-0.02	0.21	0.22	+0.01	0.90	0.90	-0.00	PDM
IF-I	186	188	+2	0.59	0.59	-0.00	0.14	0.14	+0.00	0.86	0.86	+0.00	PDM
IF-II	85	82	-3	0.84	0.83	-0.01	0.15	0.16	+0.01	0.91	0.92	+0.01	PDM
$\delta = 16/255$	clean	adv	Δ	clean	adv	Δ	clean	adv	Δ	clean	adv	Δ	
DiT-256	131	220	+89	0.37	0.26	-0.11	0.44	0.70	+0.26	0.74	0.63	-0.11	LDM
SD-V-1.4	44	225	+181	0.68	0.34	-0.34	0.22	0.68	+0.46	0.92	0.72	-0.20	LDM
SD-V-1.5	45	226	+181	0.73	0.37	-0.36	0.20	0.62	+0.42	0.94	0.78	-0.16	LDM
GD-ImageNet	109	110	+1	0.66	0.57	-0.09	0.21	0.26	+0.05	0.90	0.89	-0.01	PDM
IF-Ĭ	186	188	+2	0.59	0.58	-0.01	0.14	0.15	+0.01	0.86	0.87	+0.01	PDM
IF-II	85	86	+1	0.84	0.76	-0.08	0.15	0.21	+0.06	0.91	0.95	+0.04	PDM

Table 1: Quantitative Measurement of PGD-based Adv-Attacks for LDMs and PDMs: gradient-based diffusion attacks can attack LDMs effectively, making the difference  $\Delta$  across all evaluation metrics between edited clean image and edited adversarial image large, which means the quality of edited images drops dramatically. However, the PDMs are not affected much by the crafted adversarial perturbations, showing small  $\Delta$  before and after the attacks.

#### 3.2 EMPIRICAL ROBUSTNESS OF PDMS

To further verify the empirical robustness of pixel-space diffusion models, we proceed by designing more adaptive attacks for PDMs. We adopt some design code from (Tramer et al., 2020) to craft adaptive attacks. We first divide the attacks into two categories (C1): attack the full cascaded pipeline,

which is an end-to-end attack for the targeted editing pipeline. (C2): use diffusion loss as the objective, which follows Equation 1.

We also explore common practices for attacking randomized functions, including Expectation over Transformation (EoT) (Athalye et al., 2018), targeted attacks, and latent attacks that manipulate intermediate features. To test the robustness of Guided Diffusion (GD), we design seven combinations: textcolorblueAttack (1) and Attack (2) correspond to C1 attacks with and without EoT; Attack (3) and Attack (4) are C2 targeted and untargeted attacks; Attack (5) and Attack (6) extend the C2 targeted and untargeted attacks with EoT; and finally, Attack (7) is a latent attack (and its variant) as proposed in (Shih et al., 2024).

Attack	Description
Attack (1)	C1 attack with EoT
Attack (2)	C1 attack without EoT
Attack (3)	C2 targeted attack
Attack (4)	C2 untargeted attack
Attack (5)	C2 targeted attack with EoT
Attack (6)	C2 untargeted attack with EoT
Attack (7)	Latent attack (and its variant) (Shih et al., 2024)

Table 2: Summary of the seven adaptive attacks, see Appendix C for detailed python-format code for how (un)targeted attacks and EoT work.

Figure 3 shows that crafted perturbations fail to induce chaotic generation outcomes in PDMs. It means all the adaptive strategies widely used in the literature fail to fool PDMs. We provide numerical results of FID before and after attack in Appendix D. Moresults are shown in Appendix Figure 16. This indicates that the diffusion denoising process in the pixel space is highly robust against **empirically** existing perturbations. In the next section we will try to make explanations of the robustness of PDMs and vulnerability of LDMs.

Recent work by (Shih et al., 2024) introduces latent attacks that can effectively deceive diffusion models. The core idea is to target the intermediate layers of the U-Net architecture in Guided Diffusion (GD). While this type of attack appears capable of misleading the PDM to edit the object as something different (see Figure 4), it suffers from two major limitations: The perturbation magnitude is excessively large, with  $\ell_{\infty}(\delta) > 150/255$ . As a result, the appearance of the objects is significantly altered and further degraded by added Gaussian noise. Consequently, the diffusion model will not be able to correctly identify the object. For instance, as shown in the last block of Figure 4, when large Gaussian noise is introduced, the diffusion model mistakenly identifies the chicken as a turtle. Additionally, such latent attacks are ineffective when the editing strength is low (See Figure 16 in Appendix), indicating that the attack mechanism heavily relies on the magnitude of noise applied. In contrast, attacks against Latent Diffusion Models (LDMs) can remain effective even with small perturbation steps, as they are capable of crafting strong adversarial attacks despite limited noise being added.

#### 3.3 Denoised Diffusion tends to be Robust, the Encoder is Vulnerable

The previous two sections demonstrate that PDMs exhibit significantly stronger empirical robustness compared to LDMs. Here we provide some possible direction of explanations by looking at the loss for the attacker. First, we define the adversarial loss for attacking a diffusion model:

$$\mathcal{L}(x) = \mathbb{E}_{t,\epsilon} [f_{\theta}(x + \sigma_t \epsilon, t)],$$

which serves as a generalized form. A particular instance,  $\mathcal{L}_{adv}(x)$ , arises by choosing

$$f = \left\| \epsilon_{\theta}(x_t, t) - \epsilon \right\|_2^2$$

An important observation is that  $\epsilon$  and x appear additively in the loss term f. The key rationale behind taking the expectation over  $\epsilon$  is that, from the attacker's perspective, the exact noise  $\epsilon$  used during inference is unknown. Hence, averaging over all possible  $\epsilon$  values provides the best adversarial strategy.

The attacker wants to maximize  $\mathcal{L}(x + \delta)$  with a small perturbation  $\delta$ . However, Appendix A3 (Chen et al.) gives the bound for the Lipschitz constant of the above  $\mathcal{L}$ :

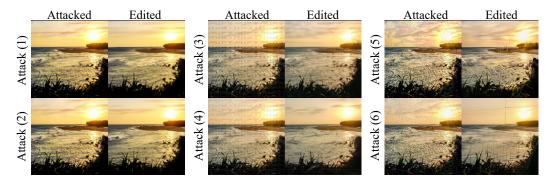


Figure 3: Crafting Adaptive Attacks for PDMs (Attack (1-6)): PDM shows robustness against end-to-end attacks and sampling based attacks, for EoT settings. We use the images in (Zheng et al., 2023) as the targeted image in the pixel space.

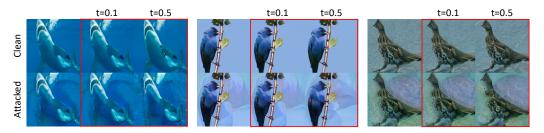


Figure 4: Latent Attacks for PDMs (Attack (7)): (Shih et al., 2024) proposes to attack the intermediate feature of the denoiser, and use a additional encoder-decoder to regularize the perturbation. This kind of attack need large perturbation  $\ell_{\infty} > 150/255$ , and it barely work for small editing steps. More results are in Appendix Figure 16.

$$\max_{\delta} \left( \mathcal{L}(x+\delta) - \mathcal{L}(x) \right) \leq \frac{1}{\sqrt{2\pi}} \|\delta\|_2 \mathbb{E}_t \left[ \frac{1}{\sigma_t} \right] \tag{4}$$

 $\frac{1}{\sqrt{2\pi}}\mathbb{E}_t[\frac{1}{\sigma_t}]$  can be roughly estimated, e.g. with EDM (Karras et al., 2022) settings of  $\sigma_t$ . The final Lipschitz constant is around 0.1. It means that the attacker suffers from the smoothness of the best guess  $\mathcal{L}$  to attack the diffusion model.

It should be noted that several studies have demonstrated successful attacks on the DiffPure pipeline (Kang et al., 2024; Xue et al., 2024; Kassis et al., 2024), primarily because the classifier itself remains vulnerable: even minor perturbations in g(x) can mislead it, which does not contradict our insights. In the context of protection attack diffusion-based mimicry, such small modifications to g(x) do not significantly degrade perceptual quality for human observers

On the other side, the vulnerability of the LDMs is caused by the vulnerability of the latent space (Xue et al., 2023), meaning that although we may set budgets for perturbations in the pixel space, the perturbations in the latent space can be large. In (Xue et al., 2023), the authors show statistics of perturbations in the latent space over the perturbations in the pixel space and this value  $\frac{|z-z'|}{|x-x'|}$  can be as large as 10, making the inputs into the denoiser ( $z_t = q_t(z), z_t' = q_t'(z')$ ) have smaller overlap (Figure 1 Middle). If we write down the distribution overlap for PDM and LDM, we have:

$$KL_{LDM}(q_t, q_t') = \frac{\|z - z'\|^2}{2\sigma^2}, \quad KL_{PDM}(q_t, q_t') = \frac{\|x - x'\|^2}{2\sigma^2}$$
 (5)

where  $\sigma^2$  is the variance at some timestep. since  $\frac{|z-z'|}{|x-x'|}$  is very large, the overlap for LDMs before and after perturbation is much smaller compared to PDMs.

Almost all the copyright protection perturbations (Shan et al., 2023; Liang and Wu, 2023; Zheng et al., 2023) are based on the insight that it is easy to craft adversarial examples to fool diffusion

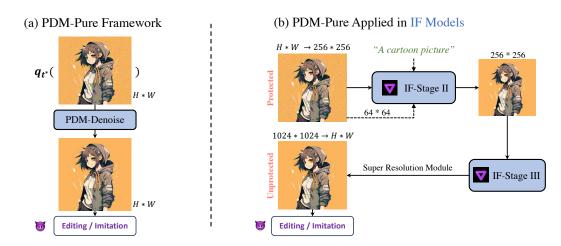


Figure 5: **PDM-Pure** is Easy to Design: (a) PDM-Pure applies SDEdit Meng et al. (2021) in the pixel space: it first runs forward diffusion with a small step  $t^*$  and then runs the denoising process. (b) We adapt the framework to DeepFloyd-IF Shonenkov et al., one of the strongest PDMs. PDM-Pure can effectively remove strong protective perturbations (e.g.  $\delta = 16/255$ ). The images we tested are sized  $512 \times 512$ .

models. We need to rethink the adversarial samples for diffusion models since there are a lot of PDMs that cannot be attacked easily. Next, we show that PDMs can be utilized to purify all adversarial patterns generated by existing methods in Section 4. This new landscape poses new challenges to ensure the security and robustness of diffusion-based copyright protection techniques.

#### 4 PDM-Pure: PDM as a Strong Universal Purifier

Since PDMs are robust to adversarial perturbations, a natural idea emerges: we can utilize PDMs as a universal purification network. This approach could potentially eliminate any adversarial patterns without knowing the nature of the attacks. We term this framework **PDM-Pure**, which is a general framework to deal with all the perturbations utilized nowadays. To fully harness the capabilities of PDM-Pure, we need to fulfill two basic requirements: (1) The perturbation adds an out-of-distribution pattern as reflected in existing works on adversarial purification/attacks using diffusion models (Nie et al., 2022; Xue et al., 2024) (2) The PDM being used is strong enough to represent  $p(x_0)$ , which can be largely determined by the dataset they are trained on.

We build PDM-Pure on the paradigm introduced in Diff-Pure Nie et al. (2022). While Diff-Pure focuses on removing adversarial noise generated for classifiers, claiming that diffusion models can diffuse and denoise these perturbations; our main point is that the diffusion process must occur in pixel space to ensure robustness. We show that LDM-based Diff-Pure fails when the noise is crafted specifically for the encoder (see Appendix Figure 17).

In the main paper, we adopt DeepFloyd-IF (Shonenkov et al.), the strongest pixel-space diffusion models nowadays as the purifier. We conduct experiments on purifying protected images sized  $512 \times 512$ . For images with a larger resolution, purifying in the resolution of  $256 \times 256$  may lose information. In Appendix K we show that PDM-Pure can also be applied to purify patches of high-resolution inputs, removing widely used protections like Glaze on artworks. More details about the how we run DeepFloyd-IF as the purification pipeline are in the Appendix I.

#### 4.1 Experimental Results of PDM-Pure

PDM-Pure is simple: we just run SDEdit to purify the protected image in the pixel space. Given our assumption that PDMs are quite robust, we can use PDMs trained on large-scale datasets as a universal black-box purifier. We follow the model pipeline introduced in Section 4 and purify images protected by various methods as shown in Table 3.

 PDM-Pure is effective: from Table 3 we can see that the purification will remove adversarial patterns for all the protection methods we tested, largely decreasing the FID score for the SDEdit task. Also, we test the protected images and purified images in more tasks including Image Inpainting (Song et al., 2020), Textual-Inversion (Gal et al., 2022), and LoRA customization (Hu et al., 2021). We show purification results for inpainting, LoRA and Text-Inversion in Figure 6. Both qualitative and quantitative results show that the purified images are no longer adversarial and can be effectively edited or imitated in different tasks without any obstruction.

Also, PDM-Pure shows SOTA results compared with previous purification methods, including some simple purifiers based on compression and filtering like Adv-Clean, crop-and-resize, JPEG Compression, NoisyUpScale (Hönig et al., 2024), and SDEdit-based methods like GrIDPure (Zhao et al., 2023), which uses patchified SDEdit with a small PDM (Dhariwal and Nichol, 2021). We also add LDM-Pure as a baseline to show that LDMs can not be used to purify the protected images. For GrIDPure, we use Guided-Diffusion trained on ImageNet to run patchified purification. All the experiments are conducted on the datasets collected in (Xue et al., 2023) under the resolution of  $512 \times 512$ . Results for higher resolutions are presented in Appendix K. We also test the ablation of timesteps used for PDM-Pure in Appendix Appendix L, from which we can see the sweet point of timesteps:  $t^*$  around 0.15 works well. We also find that PDM-Pure works better for cartoon pictures with larger plain color patches. For pictures with many details like oil paintings, it will lose some detail; however, generally the art style can still be learned well by LoRA from the attacker's perspective (e.g. Claude Monet-style in Appendix Figure 12).

#### 5 CONCLUSION AND FUTURE WORK

We show that *pixel-space diffusion models* (*PDMs*) exhibit strong empirical robustness to adversarial perturbations that reliably compromise latent diffusion models (*LDMs*), largely due to encoder-induced vulnerability in *LDMs* versus smoother pixel-space denoising in *PDMs*. Leveraging this, **PDM-Pure** uses a strong PDM to remove attack-agnostic perturbations, restoring editability and undermining current protection schemes based on adversarial noise. These results suggest that *pixel space acts as a barrier*, both resisting attacks and enabling practical purification. Future work will tighten theoretical guarantees for pixel-space robustness and explore defenses that do not rely on fragile adversarial perturbations.

Methods	AdvDM	AdvDM(-)	SDS(-)	SDS(+)	SDST	Photoguard	Mist	Mist-v2
Before Protection	166	166	166	166	166	166	166	166
After Protection	297	221	231	299	322	375	372	370
Crop-Resize	210	271	228	217	280	295	289	288
JPEG	296	222	229	297	320	359	351	348
Adv-Clean	243	201	204	244	243	266	282	270
LDM-Pure	300	251	235	300	350	385	380	375
GrIDPure	200	182	195	200	210	220	230	210
NoisyUpScale	165	173	164	166	180	191	200	199
PDM-Pure (ours)	161	170	165	159	179	175	178	170

Table 3: Quantiative Measurement of Different Purification Methods in Different Scale (FID-score): We compute the FID-score of edited purified images over the clean dataset. PDM-Pure achieves the best results on all protection methods, under strong protection with  $\delta=16$ . GrID-Pure Zhao et al. (2023) can also perform reasonably, but the performance is limited because the PDM they used is not strong enough.

# REFERENCES

N. Ahn, W. Ahn, K. Yoo, D. Kim, and S.-H. Nam. Imperceptible protection against style imitation from diffusion models. *arXiv preprint arXiv:2403.19254*, 2024.

S. Andersen. Us district court for the northern district of california. January 2023.



Figure 6: **PDM-Pure for inpainting, textual inversion and LoRA**: we show that PDM-Pure can effectively remove the adversarial perturbations of previous SOTA protection methods, undermines the effectiveness of using adversarial perturbation as protection against diffusion-based mimicry.

492

497

498

499 500

501 502

503

504

505

506 507

508

509 510

511

512

513

514

517

518

519

523

527

528

529

530

533

- A. Athalye, N. Carlini, and D. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pages 274–283. PMLR, 2018.
- 490 Y. Balaji, S. Nah, X. Huang, A. Vahdat, J. Song, K. Kreis, M. Aittala, T. Aila, S. Laine, B. Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.
- N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. In 2017 ieee symposium on security and privacy (sp), pages 39–57. Ieee, 2017.
- H. Chen, Y. Dong, S. Shao, Z. Hao, X. Yang, H. Su, and J. Zhu. Diffusion models are certifiably robust classifiers.
   In The Thirty-eighth Annual Conference on Neural Information Processing Systems.
  - J. Chen, J. Dong, and X. Xie. Exploring adversarial attacks against latent diffusion model from the perspective of adversarial transferability. arXiv preprint arXiv:2401.07087, 2024.
  - J. S. Choi, K. Lee, J. Jeong, S. Xie, J. Shin, and K. Lee. Diffusionguard: A robust defense against malicious diffusion-based image editing. *arXiv* preprint arXiv:2410.05694, 2024.
    - J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009.
    - P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
  - R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv* preprint arXiv:2208.01618, 2022.
    - I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572, 2014.
    - M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33:6840–6851, 2020.
  - J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.
- R. Hönig, J. Rando, N. Carlini, and F. Tramèr. Adversarial perturbations cannot reliably protect artists from generative ai. *arXiv* preprint arXiv:2406.12027, 2024.
  - E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Q. Huang, D. S. Park, T. Wang, T. I. Denk, A. Ly, N. Chen, Z. Zhang, Z. Zhang, J. Yu, C. Frank, et al. Noise2music: Text-conditioned music generation with diffusion models. arXiv preprint arXiv:2302.03917, 2023.
  - J. Kang, H. Yang, Y. Cai, H. Zhang, X. Xu, Y. Du, and S. He. Sita: Structurally imperceptible and transferable adversarial attacks for stylized image generation. *IEEE Transactions on Information Forensics and Security*, 2025.
- M. Kang, D. Song, and B. Li. Diffattack: Evasion attacks against diffusion-based adversarial purification.
   Advances in Neural Information Processing Systems, 36, 2024.
  - T. Karras, M. Aittala, T. Aila, and S. Laine. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35:26565–26577, 2022.
- 535
  536
  A. Kassis, U. Hengartner, and Y. Yu. Diffbreak: Is diffusion-based purification robust? arXiv preprint arXiv:2411.16598, 2024.
- N. Kumari, B. Zhang, R. Zhang, E. Shechtman, and J.-Y. Zhu. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1931–1941, 2023.

A. Li, Y. Mo, M. Li, and Y. Wang. Pid: prompt-independent data protection against latent diffusion models.
 arXiv preprint arXiv:2406.15305, 2024.

542 543

544

546

547 548

549

550

552

554 555

556

557

558

559

562

563

564

565

566 567

570

571

572

573

574

575

577 578

579

580 581

582

583

584

585

589

- C. Liang and X. Wu. Mist: Towards improved adversarial examples for diffusion models. *arXiv preprint arXiv:2305.12683*, 2023.
- C. Liang, X. Wu, Y. Hua, J. Zhang, Y. Xue, T. Song, Z. Xue, R. Ma, and H. Guan. Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples. In *International Conference on Machine Learning*, pages 20763–20786. PMLR, 2023.
- C.-H. Lin, J. Gao, L. Tang, T. Takikawa, X. Zeng, X. Huang, K. Kreis, S. Fidler, M.-Y. Liu, and T.-Y. Lin. Magic3d: High-resolution text-to-3d content creation. *arXiv* preprint arXiv:2211.10440, 2022.
- Y. Liu, C. Fan, Y. Dai, X. Chen, P. Zhou, and L. Sun. Toward robust imperceptible perturbation against unauthorized text-to-image diffusion-based synthesis. *arXiv* preprint arXiv:2311.13127, 3, 2023.
- A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
  - C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2021.
  - G. Mittal, J. Engel, C. Hawthorne, and I. Simon. Symbolic music generation with diffusion models. *arXiv* preprint arXiv:2103.16091, 2021.
- W. Nie, B. Guo, Y. Huang, C. Xiao, A. Vahdat, and A. Anandkumar. Diffusion models for adversarial purification.
   arXiv preprint arXiv:2205.07460, 2022.
  - T. C. Ozden, O. Kara, O. Akcin, K. Zaman, S. Srivastava, S. P. Chinchali, and J. M. Rehg. Optimization-free image immunization against diffusion-based editing. *arXiv preprint arXiv:2411.17957*, 2024.
  - X. Pan, P. Qin, Y. Li, H. Xue, and W. Chen. Synthesizing coherent story with auto-regressive latent diffusion models. *arXiv preprint arXiv:2211.10950*, 2022.
  - W. Peebles and S. Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
    - D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. arXiv preprint arXiv:2307.01952, 2023.
    - B. Poole, A. Jain, J. T. Barron, and B. Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *The Eleventh International Conference on Learning Representations*, 2023.
    - T. Rahman, H.-Y. Lee, J. Ren, S. Tulyakov, S. Mahajan, and L. Sigal. Make-a-story: Visual memory conditioned consistent story generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2493–2502, 2023.
  - R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
  - H. Salman, A. Khaddaj, G. Leclerc, A. Ilyas, and A. Madry. Raising the cost of malicious ai-powered image editing. arXiv preprint arXiv:2302.06588, 2023.
  - P. Sandoval-Segura, J. Geiping, and T. Goldstein. Jpeg compressed images can bypass protections against ai editing. arXiv preprint arXiv:2304.02234, 2023.
- C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis,
   M. Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models.
   Advances in Neural Information Processing Systems, 35:25278–25294, 2022.
  - R. Setty. Ai art generators hit with copyright suit over artists' images. January 2023.
- S. Shan, J. Cryan, E. Wenger, H. Zheng, R. Hanocka, and B. Y. Zhao. Glaze: Protecting artists from style mimicry by text-to-image models. *arXiv preprint arXiv:2302.04222*, 2023.
- K. Shen, R. Quan, J. Miao, J. Xiao, and Y. Yang. Tarpro: Targeted protection against malicious image editing. arXiv preprint arXiv:2503.13994, 2025.

- C.-Y. Shih, L.-X. Peng, J.-W. Liao, E. Chu, C.-F. Chou, and J.-C. Chen. Pixel is not a barrier: An effective evasion attack for pixel-domain diffusion models. *arXiv preprint arXiv:2408.11810*, 2024.
- A. Shonenkov, M. Konstantinov, D. Bakshandaeva, C. Schuhmann, K. Ivanova, and N. Klokova. IF. https://github.com/deep-floyd/IF.
- U. Singer, A. Polyak, T. Hayes, X. Yin, J. An, S. Zhang, Q. Hu, H. Yang, O. Ashual, O. Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022.
- Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. *arXiv* preprint arXiv:2011.13456, 2020.
- F. Tramer, N. Carlini, W. Brendel, and A. Madry. On adaptive attacks to adversarial example defenses. *Advances in neural information processing systems*, 33:1633–1645, 2020.
- Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- H. Xue, C. Liang, X. Wu, and Y. Chen. Toward effective protection against diffusion-based mimicry through score distillation. In *The Twelfth International Conference on Learning Representations*, 2023.
- H. Xue, A. Araujo, B. Hu, and Y. Chen. Diffusion-based adversarial sample generation for improved stealthiness and controllability. Advances in Neural Information Processing Systems, 36, 2024.
- C. Zhang, C. Zhang, M. Zhang, and I. S. Kweon. Text-to-image diffusion model in generative ai: A survey. arXiv preprint arXiv:2303.07909, 2023.
- R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.
- Z. Zhao, J. Duan, K. Xu, C. Wang, R. Z. Z. D. Q. Guo, and X. Hu. Can protective perturbation safeguard personal data from being exploited by stable diffusion? *arXiv* preprint arXiv:2312.00084, 2023.
- B. Zheng, C. Liang, X. Wu, and Y. Liu. Understanding and improving adversarial attacks on latent diffusion model. *arXiv preprint arXiv:2310.04687*, 2023.

# **Appendix**

### A BROADER IMPACT

We present significant insights in two crucial areas: adversarial machine learning research on generative diffusion models, and the protection of copyright against the malicious use of diffusion models. While existing works have revealed the vulnerability of latent diffusion models, we show that the general diffusion model in the pixel space is quite robust. PDMs reveal two new threats to the safety application of diffusion models: (1) since PDMs are robust and no existing perturbation can effectively attack them, it means that copyright protection against PDMs cannot be easily achieved with existing protective perturbations (2) PDMs can be used to purify the protective noise used to protect the LDMs, meaning that the current protection for LDMs can be bypassed. We still have a long way to go to achieve good protection against diffusion models, and more efforts should be dedicated to enhancing copyright protection for PDMs and making current protective measures more robust and reliable.

### B DETAILS ABOUT DIFFERENT DIFFUSION MODELS IN THIS PAPER

Here we introduce the diffusion models used in this work, which cover different types of diffusion (LDM, PDM), different training datasets, different resolutions, and different model structures (U-Net, Transformer):

**Guided Diffusion (PDM)** We use the implementation and checkpoint from https://github.com/openai/guided-diffusion, the Guided Diffusion models we used are trained on ImageNet (Deng et al., 2009) in resolution  $256 \times 256$ , the editing results are tested on sub-dataset of ImageNet validation set sized 500.

**IF-Stage I (PDM)** This is the first stage of the cascaded DeepFloyd IF model (Shonenkov et al.) from https://github.com/deep-floyd/IF. It is trained on LAION 1.2B with text annotation. It has a resolution of  $64 \times 64$ . the editing results are tested on the image dataset introduced in (Xue et al., 2023), including 400 anime, portrait, landscape, and artwork images.

**IF-Stage II (PDM)** This is the second stage of the cascaded DeepFloyd IF model (Shonenkov et al.) from https://github.com/deep-floyd/IF. It is a conditional diffusion model in the pixel space with  $256 \times 256$ , which is conditioned on  $64 \times 64$  low-resolution images. During the attack, we freeze the image condition and only attack the target image to be edited.

**Stable Diffusion V-1.4 (LDM)** It is one of the most popular LDMs from https://huggingface.co/CompVis/stable-diffusion-v1-4, also trained on text-image pairs, which has been widely studied in this field. It supports resolutions of  $256 \times 256$  and  $512 \times 512$ , both can be easily attacked. The encoder first encodes the image sized  $H \times W$  into the latent space sized  $4 \times H/4 \times W/4$ , and then uses U-Net combined with cross-attention to run the denoising process.

**Stable Diffusion V-1.5 (LDM)** It has the same structure as Stable Diffusion V-1.4, which is also stronger since it is trained with more steps, from https://huggingface.co/runwayml/stable-diffusion-v1-5.

**DiT-XL** (**LDM**) It is another popular latent diffusion model, that uses the backbone of the Transformer instead of the U-Net. We use the implementation from the original repository <a href="https://github.com/facebookresearch/DiT/">https://github.com/facebookresearch/DiT/</a>.

# C DETAILS ABOUT ADAPTIVE ATTACKS

```
1 # f: denoiser
2 # sdeit(f, x, t): SDEdit, t is the same during imitation
```

```
702
     3 # eps: l-inf bound for the perturbation
703
     4 # x: clean sample
     5 # delta: perturbation to be optimized
705
     6 # N: iterations
706
     8 def attack_1(): # C1
707
           delta = random_init()
708
           for _ in range(N):
    10
709
    11
               loss = sdeit(f, x + delta, t)
710
    12
               loss.backward()
               delta.update()
711
               delta.clip(epsilon)
712
           return delta
    15
713
714
    17
715 18 def attack_2(): # C1 + EoT
           delta = random_init()
    19
716
           for _ in range(N):
    20
717
    21
               grad = 0
718
               for _ in EoT_round_num:
    22
719 23
                   loss = sdeit(f, x + delta, t)
                   loss.backward()
720 24
                   grad.update()
    25
721
               delta.update()
722
               delta.clip()
    27
723
          return delta
    28
724
725
    31 def attack_3(): # C2 + untargeted
726
           delta = random_init()
    32
727
           for _ in range(N):
    33
728
    34
               loss = -MSE(eps - f(x + eps, t))
729
    35
              loss.backward()
              delta.update()
730
    36
               delta.clip()
    37
731
    38
           return delta
732
    39
733
    40
734 41 def attack_4(): # C2 + targeted
          delta = random_init()
735 42
           z = random_noise()
736
    44
          for _ in range(N):
737
               loss = MSE(f(x + eps, t) - z)
    45
738
               loss.backward()
    46
739
               delta.update()
740
    48
               delta.clip()
    49 return delta
741
```

# D NUMERICAL RESULT OF FID FOR ADAPTIVE ATTACKS AGAINST PDM

$\delta = 8$	Clean	Attack1	Attack2	Attack3	Attack4	Attack5	Attack6
FID	109	109	110	111	112	111	110
$\delta = 16$	Clean	Attack1	Attack2	Attack3	Attack4	Attack5	Attack6
FID	109	109	111	110	111	109	110

# E DETAILS ABOUT DIFFERENT PROTECTION METHODS IN THIS PAPER

We introduce different protection methods tested in this paper, of which all the original versions are designed for LDMs. All the adversarial attacks work under white box settings of PGD-attack, varying from each other with different adversarial losses:

**AdvDM** AdvDM is one of the first adversarial attacks proposed in (Liang et al., 2023), it used a Monte-Carlo-based adversarial loss which can effectively attack latent diffusion models, we also call this loss semantic loss:

$$\mathcal{L}_S(x) = \mathbb{E}_{t,\epsilon} \mathbb{E}_{z_t \sim q_t(\mathcal{E}_{\phi}(x))} \| \epsilon_{\theta}(z_t, t) - \epsilon \|_2^2$$
(6)

**PhotoGuard** PhotoGuard is proposed in (Salman et al., 2023), it takes the encoder, making the encoded image close to a target image y, we also call it textural loss:

$$\mathcal{L}_T(x) = -\|\mathcal{E}_{\phi}(x) - \mathcal{E}_{\phi}(y)\|_2^2 \tag{7}$$

**Mist** Mist (Liang and Wu, 2023) finds that  $L_T(x)$  can better enhance the attacks if the target image y is chosen to be periodical patterns, the final loss combined  $L_T(x)$  and  $L_S(x)$ :

$$\mathcal{L} = \lambda L_T(x) + L_S(x) \tag{8}$$

**SDS(+)** Proposed in (Xue et al., 2023), it is proven to be a more effective attack compared to the original AdvDM, where the gradient  $\nabla_x \mathcal{L}(x)$  is expensive to compute. By using the score distillation-based loss, it shows good performance and remains effective at the same time:

$$\nabla_x \mathcal{L}_{SDS}(x) = \mathbb{E}_{t,\epsilon} \mathbb{E}_{z_t} \left[ \lambda(t) (\epsilon_{\theta}(z_t, t) - \epsilon) \frac{\partial z_t}{\partial x_t} \right]$$
 (9)

**SDS(-)** Similar to SDS(+), it swaps gradient ascent in the original PGD with gradient descent, which turns out to be even more effective.

$$\nabla_x \mathcal{L}_{SDS(-)}(x) = -\mathbb{E}_{t,\epsilon} \mathbb{E}_{z_t} \left[ \lambda(t) (\epsilon_\theta(z_t, t) - \epsilon) \frac{\partial z_t}{\partial x_t} \right]$$
 (10)

**Mist-v2** It was proposed in (Zheng et al., 2023) using the Improved Targeted Attack (ITA), which turns out to be very effective, especially when the budget is small. It is also more effective to attack LoRA:

$$\mathcal{L}_{S}(x) = \mathbb{E}_{t,\epsilon} \mathbb{E}_{z_{t} \sim q_{t}(\mathcal{E}_{\phi}(x))} \|\epsilon_{\theta}(z_{t}, t) - z_{0}\|_{2}^{2}$$

$$\tag{11}$$

where  $z_0 = \mathcal{E}(y)$  is the latent of the target image, which is the same as the typical image used in Mist.

Glaze It is the most popular protection claimed to safeguard artists from unauthorized imitation (Shan et al., 2023) and is widely used by the community. while it is not open-sourced, it also attacks the encoder like the Photoguard. Here we only test it in the purification stage, where we show that the protection can also be bypassed.

**End-to-End Attack** It is also first proposed in (Salman et al., 2023), which attacks the editing pipeline in a end-to-end manner. Although it is strong, it is not practical to use and does not show dominant privilege compared with other protection methods.

# F DETAILS ABOUT THE LATENT ATTACKS FOR PDMS

In an attempt to extend the latent-space attacks onto PDMs, (Shih et al., 2024) introduces atkPDM+. This method uses a pre-trained VAE to attack the PDM by extracting feature vectors from the encoder network. The attack optimizes the latent vector with a Wasserstein distance objective calculated at the VAE middle layer activations:

$$\mathcal{L}_{attack}(x_t, x_t^{adv}) = -\mathcal{W}_2(\mathcal{U}_{\theta}^{(mid)}(x_t), \mathcal{U}_{\theta}^{(mid)}(x_t^{adv}))$$

A second optimization cycle is then run to limit the change in pixel-space by optimizing the distance between the feature vector generated by a pre-trained image classifier taken from the original image and the decoded attacked latent.

We observe, however, that in this attack the perturbation is clearly visible, and the pixel-wise distance is large:  $||x - x_{adv}|| \ge 150$ .

# G DETAILS ABOUT THE EVALUATION METRICS

Here we introduce the quantitative measurement we used in our experiments:

- We measure the SDEdit results after the adversarial attacks using Fréchet Inception Distance (FID) (Heusel et al., 2017) over the relevant datasets (for models trained on ImageNet such as GD (Dhariwal and Nichol, 2021) and DiT (Peebles and Xie, 2023) we use a sub-dataset of ImageNet as the relevant dataset, for those trained on LAION, we use the collected dataset to calculate the FID). We also use Image-Alignment Score (IA-score) (Kumari et al., 2023), which can be used to calculate the cosine-similarity between the CLIP embedding of the edited image and the original image. Also, we use some basic evaluations, where we calculate the Structural Similarity (SSIM) (Wang et al., 2004) and Perceptual Similarity (LPIPS) (Zhang et al., 2018) compared with the original images.
- To measure the purification results, we test the Fréchet Inception Distance (FID) (Heusel et al., 2017) over the collected dataset compared with the dataset generated by running SDEdit over the purified images in the strength of 0.3.

#### H DETAILS ABOUT DIFFERENT PURIFICATION METHODS

**Adv-Clean:** https://github.com/lllyasviel/AdverseCleaner, a training-free filter-based method that can remove adversarial noise for a diffusion model, it works well to remove high-frequency noise.

**Crop** & **Resize:** first crops the image by 20% and then resizes the image to the original size, it turns out to be one of the most effective defense methods (Liang and Wu, 2023).

**JPEG compression:** (Sandoval-Segura et al., 2023) reveals that JPEG compression can be a good purification method, and we adopt the 65% as the quality of compression in (Sandoval-Segura et al., 2023).

**LDM-Pure:** We also try to use LDMs to run SDEdit as a naive purifier, sadly it does not work, because the adversarial protection transfers well between different LDMs.

**GrIDPure:** It is proposed in (Zhao et al., 2023) as a purifier, GrIDPure first divides an image into patches sized  $128 \times 128$ , and then purifies the 9 patches sized  $256 \times 256$ . Also, it combined the four corners sized  $128 \times 128$  to purify it so we have 10 patches to purify in total. After running SDEdit with a small noise (set to 0.1T), we reassemble the patches into the original size, pixel values are assigned using the average values of the patches they belong to. More details can be seen in (Zhao et al., 2023).

# I DETAILS ABOUT PDM-PURE

Here, we explain in detail how to adapt DeepFloyd-IF (Shonenkov et al.), the strongest open-source PDM as far as we know, for PDM-Pure. DeepFloyd-IF is a cascaded text-to-image diffusion model trained on 1.2B text-image pairs from LAION dataset (Schuhmann et al., 2022). It contains three stages named IF-Stage I, II, and III. Here we only use Stage II and III since Stage I works in a resolution of 64 which is too low. Given a perturbed image  $x_{W\times H}$  sized  $W\times H$ , we first resize it into  $x_{64\times 64}$  and  $x_{256\times 256}$ . Then we use a general prompt  $\mathcal P$  to do SDEdit (Meng et al., 2021) using the Stage II model:

 $x_t = \mathbf{IF-II}(x_{t+1}, x_{64 \times 64}, \mathcal{P}) \tag{12}$ 

where  $t=T_{\rm edit}-1,...,1,0$ ,  $x_{T_{\rm edit}}=x_{256\times256}$ . A larger  $T_{\rm edit}$  may be used for larger noise.  $x_0$  is the purified image we get in the  $256\times256$  resolution space, where the adversarial patterns should be already purified. We can then use IF Stage III to further up-sample it into  $1024\times1024$  with  $x_{1024\times1024}=$  **IF-III** $(x_0,p)$ . Finally, we can sample into  $H\times W$  as we want through downsampling. This whole process is demonstrated in Figure 5. After purification, the image is no longer adversarial to the targeted diffusion models and can be effectively used in downstream tasks.

### J MORE EXPERIMENTAL RESULTS

In this section, we present more experimental results.

#### J.1 More Visualizations of Attacking PDMs

We show more results of attacking LDMs and PDMs in Figure 7, where we attack them with a different budget  $\delta=4,8,16$ . We can see that all the LDMs can be easily attacked, while the PDMs cannot be attacked, even the largest perturbations will not fool the editing process. In fact, the editing process is trying to purify the strange perturbations.

# J.2 More Visualizations of PDM-Pure and Baseline Methods

We show more qualitative results of the proposed PDM-Pure based on IF. First, we show purified samples of PDM-Pure in Figure. 9, from which we can see that PDM-Pure can remove large protective perturbations and largely preserve details.

Compared with GrIDPure (Zhao et al., 2023), we find that PDM-Pure shows better results when the noise is large and colorful, as is illustrated in Figure 10. Also, though GrIDPure merges patches, it still shows boundary lines between patches.

Compared with other baseline purification methods such as Adv-Clean, Crop-and-Resize, and JPEG compression, PDM-Pure shows much better results (Figure 8) for different kinds of protective noise, showing that it is capable to serve as a universal purifier. We choose AdvDM, Mist, and SDS as the representative of three kinds of protection.

#### J.3 More Visualizations of PDM-Pure for Downstream Tasks

After applying PDM-Pure to the protected images, they are no longer adversarial to LDMs and can be easily edited or imitated. Here we will demonstrate more results on editing the purified images on downstream tasks.

In Figure 11, we show more results to prove that the purified images can be edited easily, and the quality of the editing results is high. It means that PDM-Pure can bypass the protection very well for inpainting tasks.

In Figure 12 we show more results on purifying Mist (Liang and Wu, 2023) and Glaze (Shan et al., 2023) perturbations, and then running LoRA customized generation. From the figure, we can see that PDM-Pure can make the protected images easy to imitate again.

#### K PDM-Pure For Higher Resolution

In this paper, we mainly apply PDM-Pure for images sized  $512 \times 512$ , which is also the most widely used resolution for latent diffusion models. When the resolution is  $512 \times 512$ , running SDEdit using Stage II of DeepFloyd makes sense, while if the image size becomes larger, details may be lost because of the downsampling. Hopefully, we can still do purification patch-by-patch with PDM-Pure, in Figure 13 we show purification results on images with different resolutions protected by Glaze (Shan et al., 2023).



Figure 7: PDMs cannot be Attacked as LDMs: we conduct experiments on various models with various budgets, even the largest budget will not affect the PDMs, showing that PDMs are adversarially robust. For each block, the first column is the attacked image, and the second and third columns are edited images, where the third column adopts larger editing strength.



Figure 8: PDM-Pure Compared With Other Baseline Methods: we test all the baselines on three typical kinds of protection methods, with  $\delta=16/255$ . PDM-Pure shows strong performance.

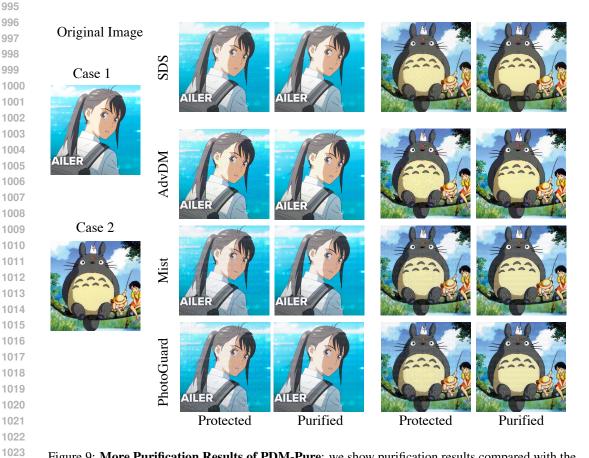


Figure 9: More Purification Results of PDM-Pure: we show purification results compared with the clean image, working on SDS, AdvDM, Mist, and PhotoGuard.

Protected

Edit (GrIDPure)

# Purification Results: PDM-Pure (IF) vs GrIDPure Protected by AdvDM Protected by Mist Protected Image Protected Image PDM-Pure (IF) Clean Image GrIDPure PDM-Pure (IF) Clean Image GrIDPure SDEdit after Purification: PDM-Pure (IF) vs GrIDPure Protected Edit (GrIDPure) Edit (PDM-Pure) Protected Edit (GrIDPure) Edit (PDM-Pure)

Figure 10: **PDM-Pure vs GrIDPure**: PDM-Pure is better than GrIDPure, especially when the adversarial pattern is strong such as AdvDM. The bottom half of this figure shows the editing results of purified images, we can see that the editing results of GrIDPure still have some artifacts.

Protected

Edit (GrIDPure)

Edit (PDM-Pure)

Edit (PDM-Pure)

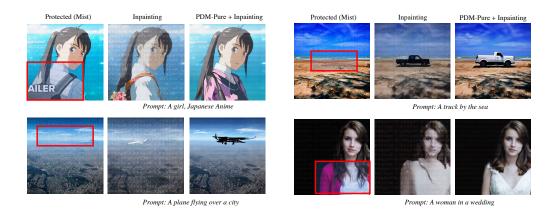


Figure 11: More Results of PDM-Pure Bypassing Protection for Inpainting: after purification, the protected images can be easily inpainted with high quality. The protective perturbations are generated using Mist with  $\delta = 16/255$ , which is a strong perturbation.

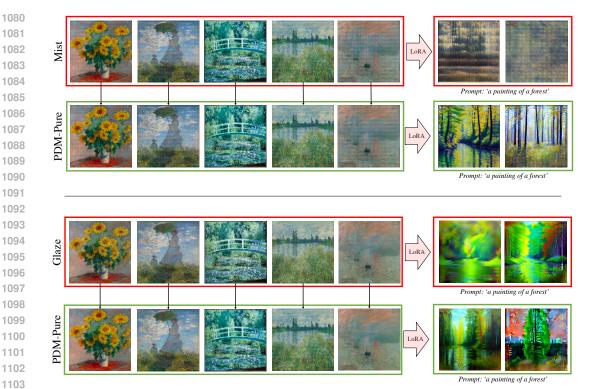


Figure 12: **More Results of PDM-Pure Bypassing Protection for LoRA**: after purification, the protected images can be imitated again. Here we show examples using 5 paintings of Claude Monet.

# L Ablations of $t^*$ in PDM-Pure

The PDM-Pure on DeepFloyd-IF we used in this paper uses the default settings of SDEdit with  $t^* = 0.1T$ . And we respace the diffusion model into 100 steps, so we only need to run 10 denoising steps. It can be run on one A6000 GPU, occupying 22G VRAM in 30 seconds.

Here we show some ablation about the choice of  $t^*$ . In fact, in many SDEdit papers,  $t^*$  can be roughly defined by trying different  $t^*$  that can be used to purify different levels of noise. We try  $t^* = 0.01, 0.1, 0.2$ , in Figure 14 we can see that when  $t^* = 0.01$  the noise is not fully purified, and when  $t^* = 0.2$ , the details in the painting are blurred. It should be noted that the sweet spot for different images and different noises can be slightly different, so one is advised to do some trials before purification.

Methods	AdvDM	AdvDM(-)	SDS(-)	SDS(+)	SDST	Photoguard	Mist	Mist-v2
Clean Attacked	0.95 0.73	0.95 0.70	0.95 0.68	0.95 0.76	0.95 0.61	0.95 0.61	0.95 0.62	0.95 0.63
PDM-Pure	0.94	0.93	0.92	0.93	0.93	0.94	0.93	0.93

Table 4: IA Score of SDEdit results After Purification

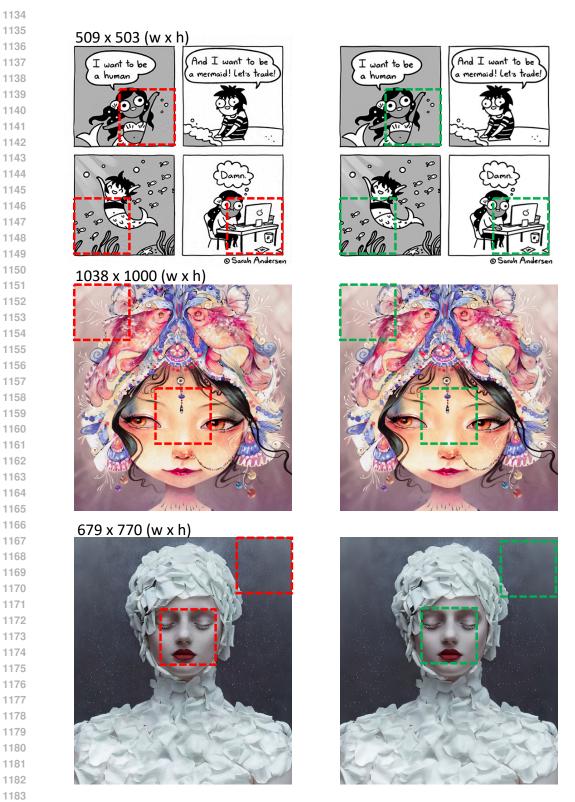


Figure 13: **PDM-Pure Working On Images with Higher Resolution**: we show the results of applying PDM-Pure for images with higher resolutions, the images are protected using Glaze (Shan et al., 2023). We can see from the figure that the adversarial patterns (in the red box) can be effectively purified (in the green box). Zoom in on the computer for a better view.

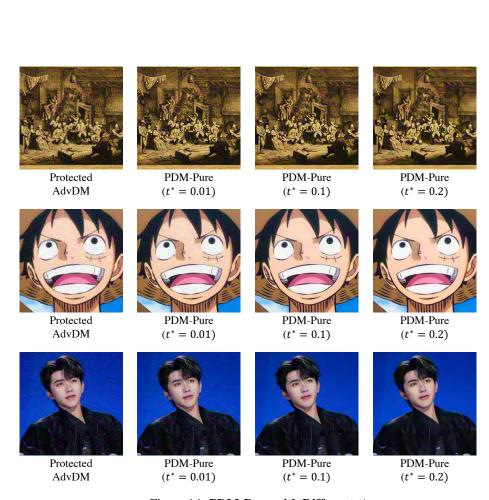


Figure 14: **PDM-Pure with Different**  $t^*$ 

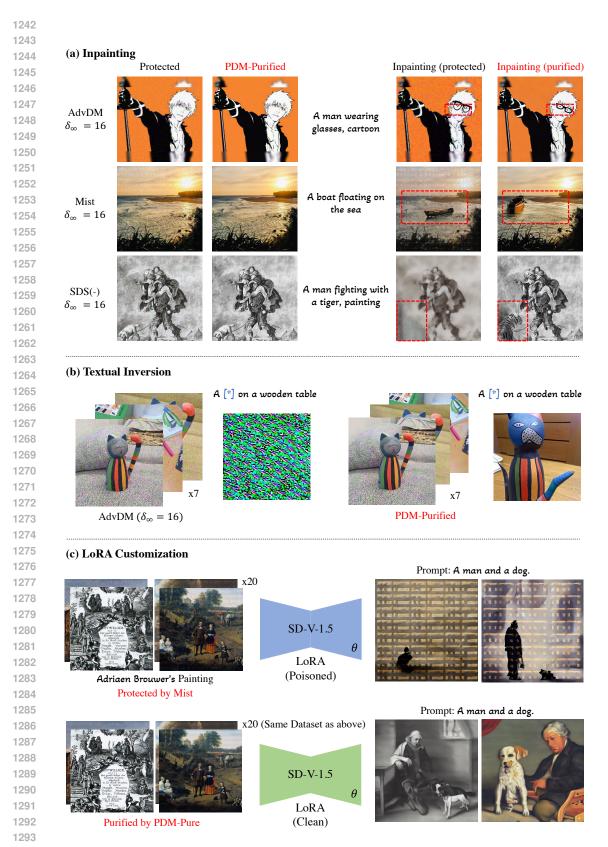


Figure 15: PDM-Pure for inpainting, textual inversion and LoRA

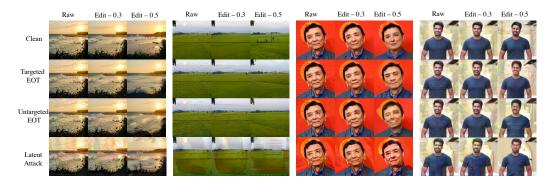


Figure 16: More results for adaptive attacks for PDM: here we show attacking results for one PDM (Guided-Diffusion (Dhariwal and Nichol, 2021)), we conduct SDEdit with two different strengths 0.3 and 0.5 to test the attacking performance. We show results for targeted/untargeted attack with gradent aggregation (Targeted/Untargeted EOT), we also show results for latent attacks following the settings in (Shih et al., 2024). We can see all the attacks is not that successful for the pixel-space diffusion model.

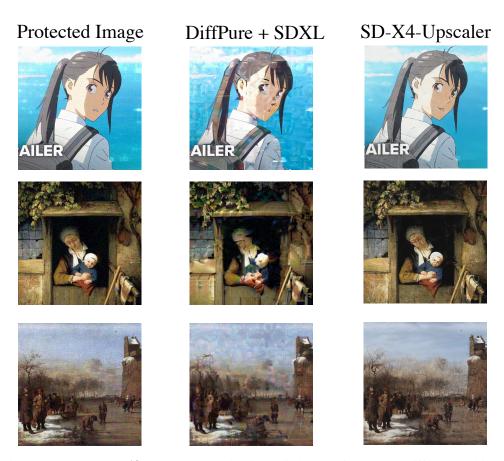


Figure 17: **LDM as Purifier**: When protection is applied to the given LDM, DiffPure combined with the LDM will fail to function effectively, as the purification process can be easily fooled. Additionally, the LDM-based upscaler lacks stability, often resulting in poor detail quality.