



AAAI-25 / IAAI-25 / EAAI-25  
FEBRUARY 25 – MARCH 4, 2025 | PHILADELPHIA, USA

# Pixel Is Not A Barrier: An Effective Evasion Attack for Pixel-Domain Diffusion Models



Chun-Yen Shih<sup>1,3,\*</sup>

Li-Xuan Peng<sup>3,\*</sup>

Jia-Wei Liao<sup>1,3</sup>

Ernie Chu<sup>2,3</sup>

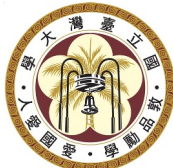
Cheng-Fu Chou<sup>1</sup>

Jun-Cheng Chen<sup>3</sup>

<sup>1</sup> National Taiwan University,

<sup>2</sup> Johns Hopkins University,

<sup>3</sup> Research Center for Information Technology Innovation, Academia Sinica



Project page



# Background

Diffusion Models allows users to generate photorealistic image with ease.

## Stable Diffusion



## ControlNet

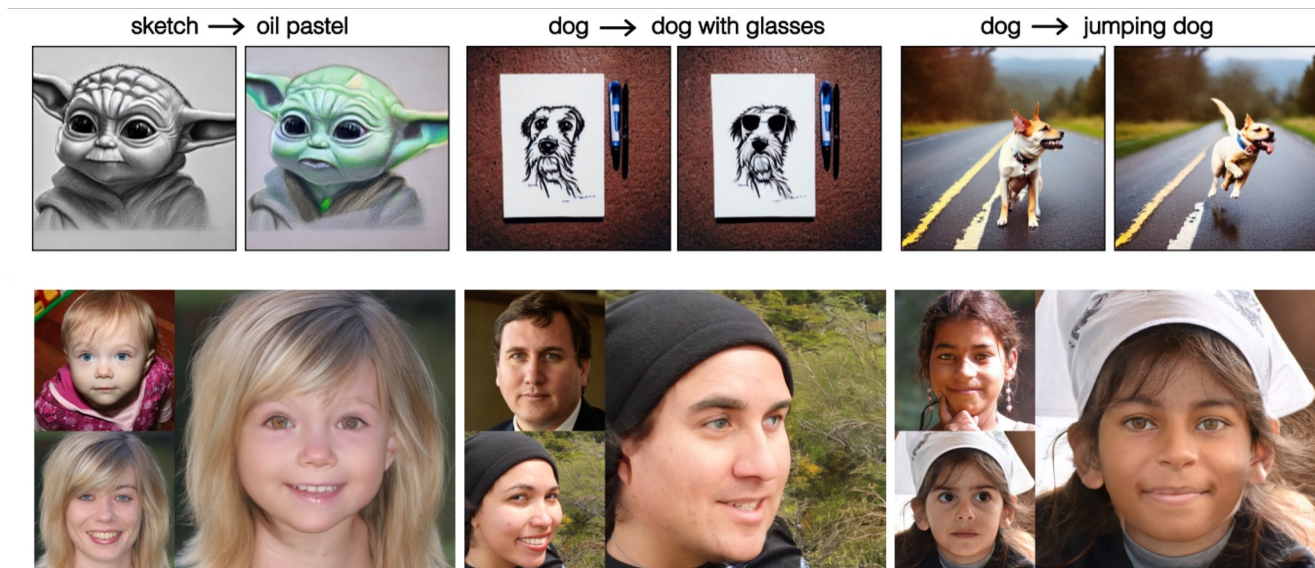
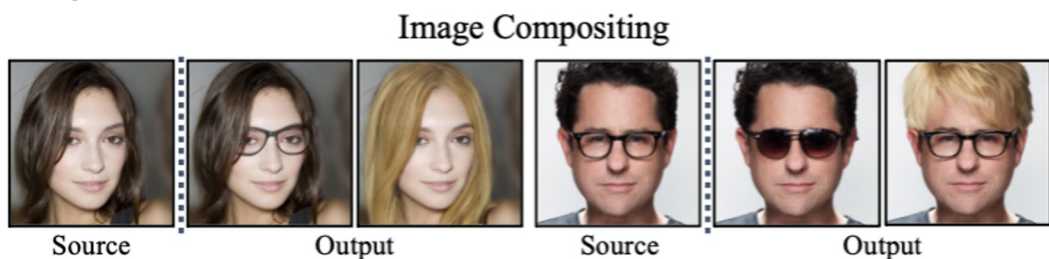
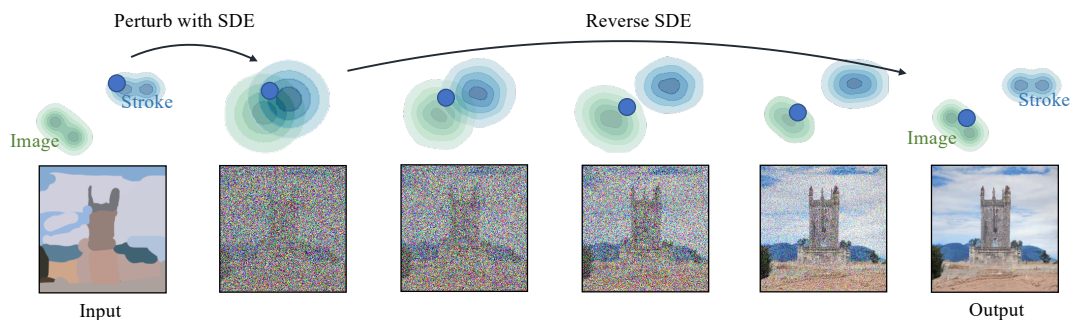


1. Robin Rombach et al. High-resolution image synthesis with latent diffusion models. CVPR 2022.
2. Lvmin Zhang et al. Adding conditional control to text-to-image diffusion models. ICCV 2023



# Background

Diffusion Models also allow easily converting image to noisy latent for image translations or editing.

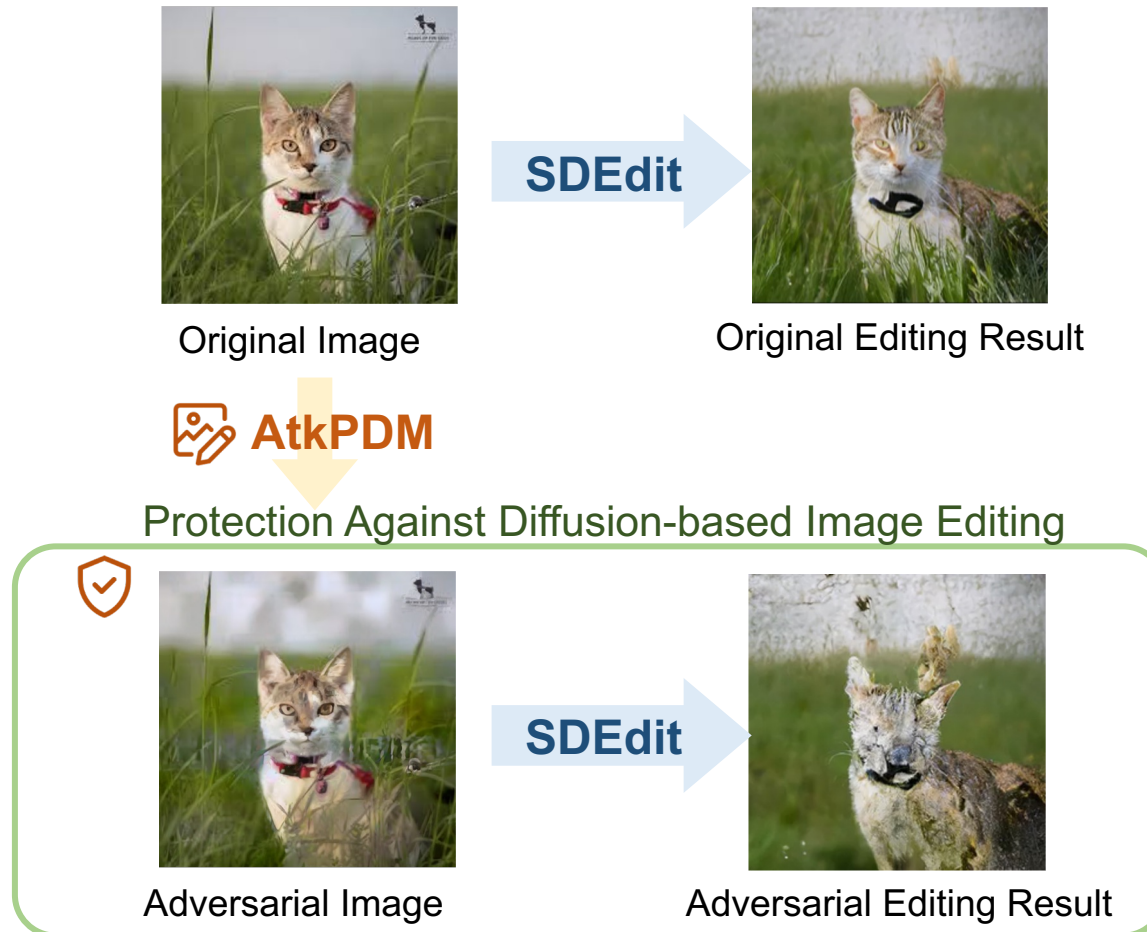


1. Chen-Lin Meng et al. SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations. ICLR 2022.
2. Gaurav Parmar et al. Zero-shot image-to-image translation. ACM SIGGRAPH 2023.
3. Wen-Liang Zhao et al. Diffswap: High-fidelity and controllable face swapping via 3d-aware masked diffusion. CVPR 2023

# Motivation of Attacking as Protection

How to protect our image against diffusion-based editing?

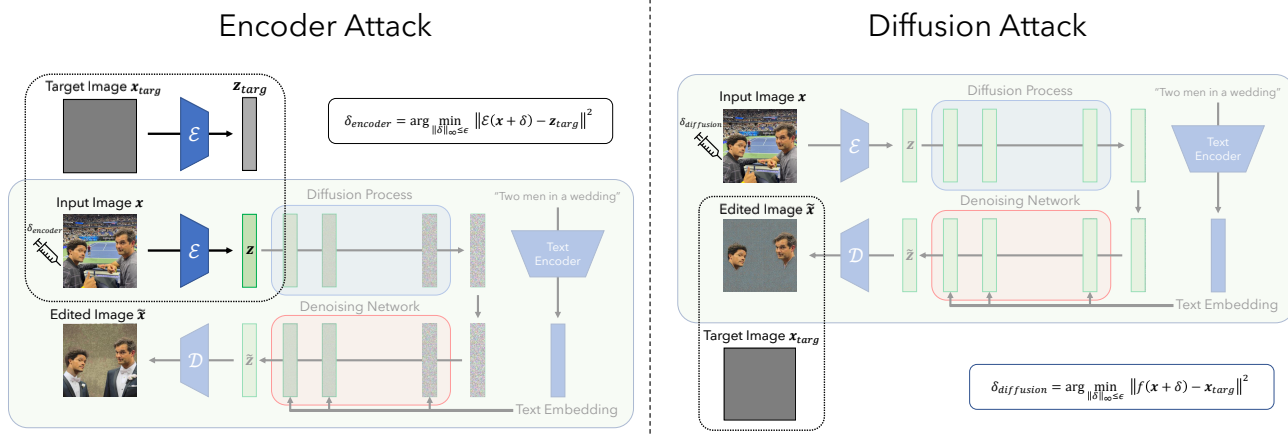
We can approach this goal as an adversarial attack to the diffusion models





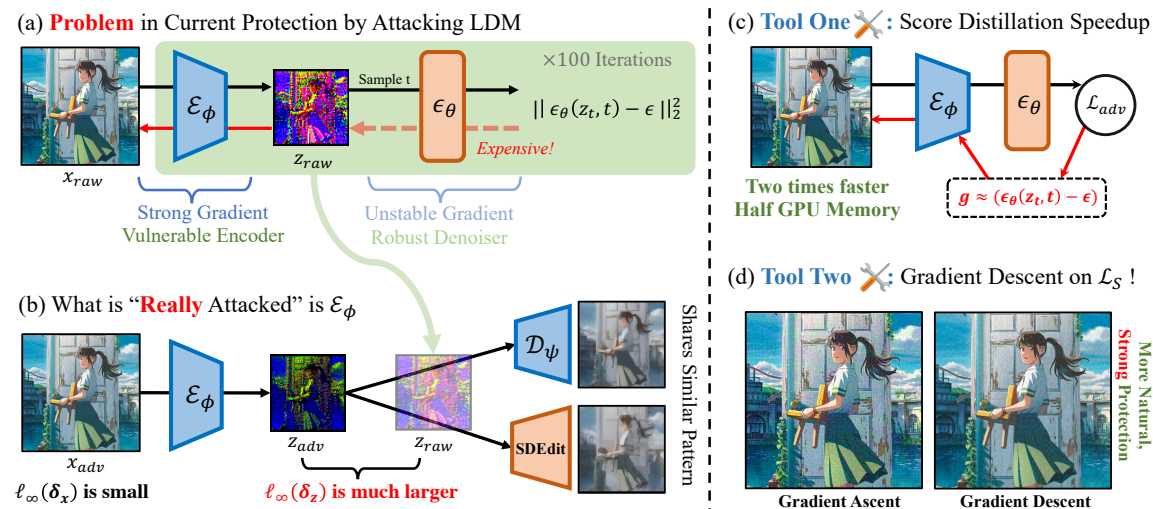
# Previous Works

## PhotoGurad [ICML 2023]



Attacking diffusion process as a whole with back-propagation requires substantial memory usage.

## Diff-Protect [ICLR 2024]

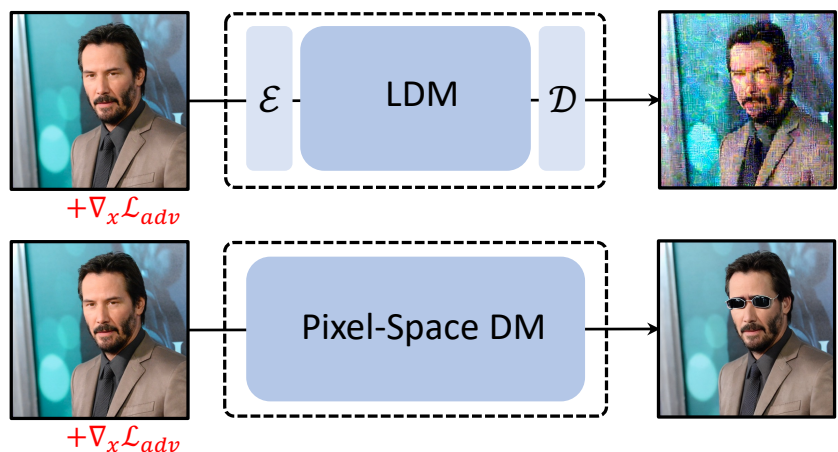


The attack effectiveness is mainly attributed to the vulnerability of the VAE encoders in LDM.

1. Hadi Salman et al. Raising the cost of malicious AI-powered image editing. ICML 2023.
2. Haotian Hue et al. Toward effective protection against diffusion-based mimicry through score distillation.. ICLR 2024

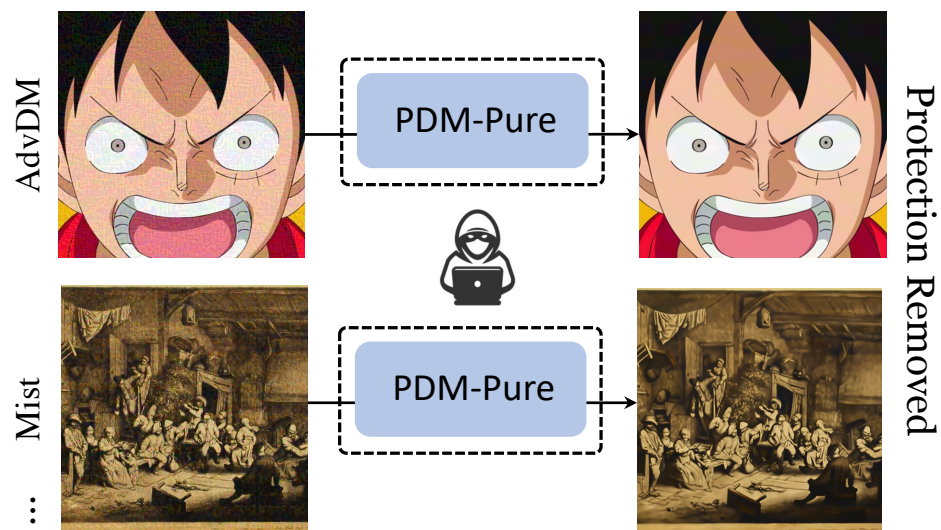
# Previous Works

(a) Adv-samples for PDMs are largely overlooked

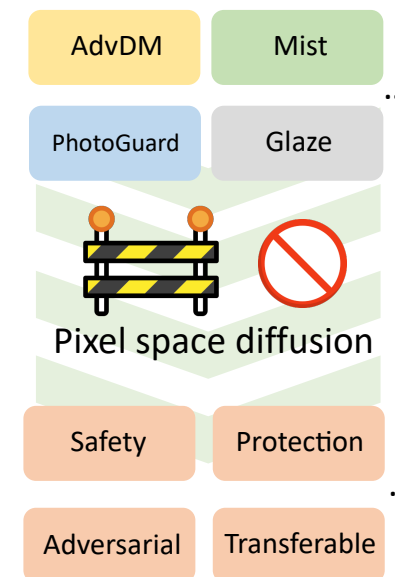


**Rethink:** LDM is easy to attack, BUT can we attack PDMs? 🤔

(b) Protections can be easily bypassed using PDM

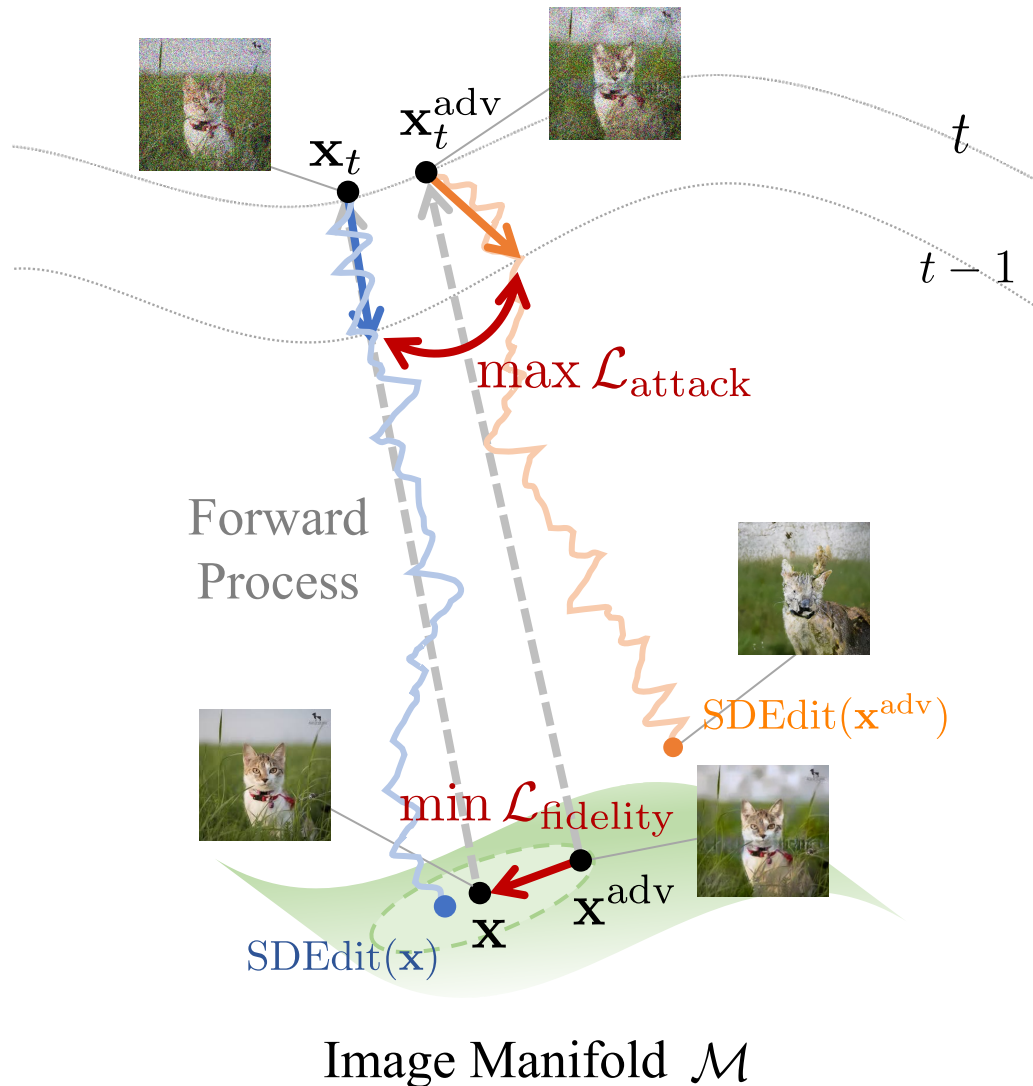


(c) Pixel is a Barrier



**Question:** Can we design an effective attack on the diffusion process that applies universally to both Pixel-based Diffusion Models (PDMs) and LDMs without relying on the vulnerability of the VAE encoder (specific to LDMs) or requiring the computational cost of back-propagating through every diffusion step?

# Problem Formulation and Methodology



## Problem

$$\max_{\mathbf{x}^{\text{adv}} \in \mathcal{M}} d(\text{SDEdit}(\mathbf{x}, t), \text{SDEdit}(\mathbf{x}^{\text{adv}}, t))$$

subject to  $d'(\mathbf{x}, \mathbf{x}^{\text{adv}}) \leq \delta$

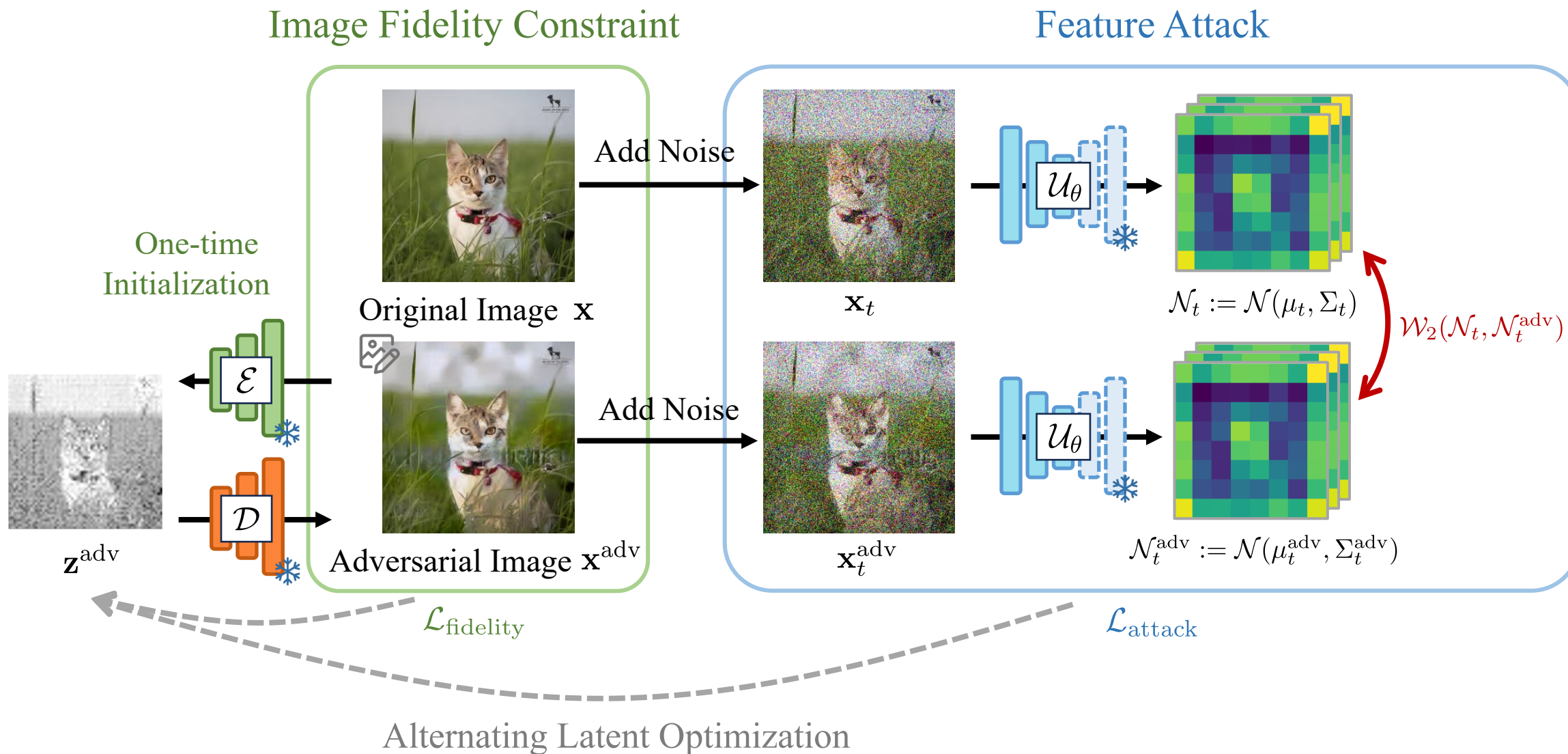
## Proposed Losses

$$\max_{\mathbf{x}^{\text{adv}} \in \mathcal{M}} \mathbb{E}_{t, \mathbf{x}_t | \mathbf{x}, \mathbf{x}_t^{\text{adv}} | \mathbf{x}} \mathcal{L}_{\text{attack}}(\mathbf{x}_t, \mathbf{x}_t^{\text{adv}})$$

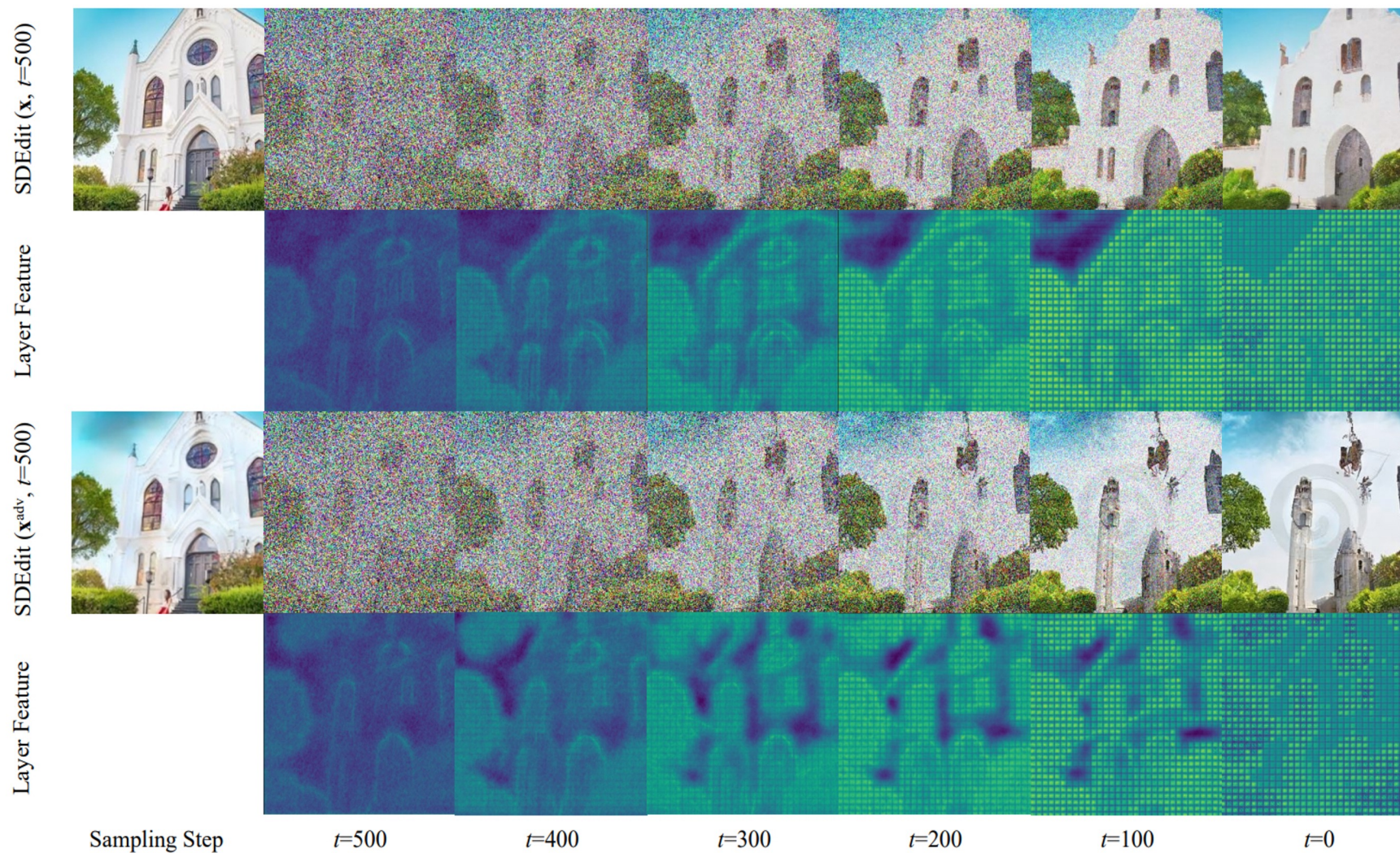
subject to  $\mathcal{L}_{\text{fidelity}}(\mathbf{x}, \mathbf{x}^{\text{adv}}) \leq \delta$



# Proposed Method



# Feature Attack Visualization





# Qualitative Results





# Quantitative Comparisons

Methods		Adversarial Image Quality			Attacking Effectiveness			
		SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\downarrow$	PSNR $\downarrow$	LPIPS $\uparrow$	IA-Score $\downarrow$
Church	AdvDM (Liang et al. 2023)	0.37 $\pm$ 0.09	28.17 $\pm$ 0.22	0.73 $\pm$ 0.16	0.89 $\pm$ 0.05	31.06 $\pm$ 1.94	0.17 $\pm$ 0.09	0.93 $\pm$ 0.04
	Diff-Protect (Xue et al. 2023)	0.39 $\pm$ 0.07	28.03 $\pm$ 0.12	0.67 $\pm$ 0.11	0.82 $\pm$ 0.05	31.90 $\pm$ 1.08	0.23 $\pm$ 0.07	0.91 $\pm$ 0.04
	AtkPDM (Ours)	<u>0.75</u> $\pm$ 0.03	<u>28.22</u> $\pm$ 0.10	<u>0.26</u> $\pm$ 0.04	<b>0.75</b> $\pm$ 0.04	<b>29.61</b> $\pm$ 0.23	<b>0.40</b> $\pm$ 0.05	<b>0.76</b> $\pm$ 0.06
	AtkPDM <sup>+</sup> (Ours)	<b>0.81</b> $\pm$ 0.03	<b>28.64</b> $\pm$ 0.19	<b>0.13</b> $\pm$ 0.02	0.79 $\pm$ 0.04	30.05 $\pm$ 0.47	0.33 $\pm$ 0.07	0.81 $\pm$ 0.06
Cat	AdvDM (Liang et al. 2023)	0.48 $\pm$ 0.09	28.34 $\pm$ 0.18	0.65 $\pm$ 0.12	0.96 $\pm$ 0.02	32.32 $\pm$ 2.49	0.10 $\pm$ 0.05	0.97 $\pm$ 0.03
	Diff-Protect (Xue et al. 2023)	0.33 $\pm$ 0.10	28.03 $\pm$ 0.15	0.80 $\pm$ 0.15	0.90 $\pm$ 0.05	33.94 $\pm$ 1.93	0.18 $\pm$ 0.08	0.95 $\pm$ 0.03
	AtkPDM (Ours)	<u>0.71</u> $\pm$ 0.06	<u>28.47</u> $\pm$ 0.18	<u>0.29</u> $\pm$ 0.05	<b>0.83</b> $\pm$ 0.03	<b>30.73</b> $\pm$ 0.51	<b>0.39</b> $\pm$ 0.05	<b>0.81</b> $\pm$ 0.04
	AtkPDM <sup>+</sup> (Ours)	<b>0.83</b> $\pm$ 0.04	<b>29.41</b> $\pm$ 0.37	<b>0.09</b> $\pm$ 0.02	0.93 $\pm$ 0.01	33.02 $\pm$ 0.74	0.18 $\pm$ 0.02	0.92 $\pm$ 0.01
Face	AdvDM (Liang et al. 2023)	0.48 $\pm$ 0.05	<b>28.75</b> $\pm$ 0.18	0.64 $\pm$ 0.10	0.99 $\pm$ 0.00	37.96 $\pm$ 1.75	0.02 $\pm$ 0.01	0.99 $\pm$ 0.00
	Diff-Protect (Xue et al. 2023)	0.25 $\pm$ 0.04	28.09 $\pm$ 0.20	0.91 $\pm$ 0.11	0.95 $\pm$ 0.02	35.33 $\pm$ 1.62	0.08 $\pm$ 0.04	0.96 $\pm$ 0.02
	AtkPDM (Ours)	<u>0.56</u> $\pm$ 0.04	28.01 $\pm$ 0.22	<u>0.36</u> $\pm$ 0.04	<b>0.74</b> $\pm$ 0.03	<b>29.14</b> $\pm$ 0.36	<b>0.40</b> $\pm$ 0.05	<b>0.62</b> $\pm$ 0.07
	AtkPDM <sup>+</sup> (Ours)	<b>0.81</b> $\pm$ 0.04	28.39 $\pm$ 0.20	<b>0.12</b> $\pm$ 0.03	0.86 $\pm$ 0.03	30.26 $\pm$ 0.72	0.24 $\pm$ 0.07	0.80 $\pm$ 0.08

Table 1: Quantitative results in attacking different unconditional PDMs. The best is marked in bold and the second best is underlined. Errors denote one standard deviation of all images in our test datasets.

Methods		Adversarial Image Quality			Attacking Effectiveness			
		SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\downarrow$	PSNR $\downarrow$	LPIPS $\uparrow$	IA-Score $\downarrow$
Diff-Protect (Xue et al. 2023)		0.47 $\pm$ 0.08	27.96 $\pm$ 0.08	0.46 $\pm$ 0.05	<b>0.49</b> $\pm$ 0.10	<b>28.13</b> $\pm$ 0.15	<b>0.36</b> $\pm$ 0.10	<b>0.79</b> $\pm$ 0.06
AtkPDM <sup>+</sup> (Ours)		<b>0.79</b> $\pm$ 0.06	<b>28.48</b> $\pm$ 0.33	<b>0.06</b> $\pm$ 0.02	0.72 $\pm$ 0.10	28.50 $\pm$ 0.48	0.10 $\pm$ 0.04	0.86 $\pm$ 0.08

Table 2: Quantitative results in attacking conditional PDM DeepFloyd IF. The best is marked in bold and the second best is underlined. Errors denote one standard deviation of all images in our test datasets.

# Quantitative Results on Defense Method and Attack Transferability

Defense Method	Attacking Effectiveness			
	SSIM ↓	PSNR ↓	LPIPS ↑	IA-Score ↓
LDM-Pure	0.78	29.84	0.35	0.80
Crop-and-Resize	0.68	29.28	0.42	0.79
JPEG Comp.	0.78	29.82	0.36	0.79
None	0.79	30.05	0.33	0.81

Table 3: Quantitative results of our adversarial images against defense methods. LDM-Pure, Crop-and-Resize, and JPEG Compression fail to defend our attack. “None” indicates no defense is applied, as the baseline for comparison.

Setting	Attacking Effectiveness			
	SSIM ↓	PSNR ↓	LPIPS ↑	IA-Score ↓
White Box	0.79	30.05	0.33	0.81
Black Box	0.86	30.25	0.29	0.85
Difference	0.07	0.20	0.04	0.04

Table 4: Quantitative results of black box attack. We use the same set of adversarial images and feed to white box and black box models to examine the black box transferability.



# Ablation Study



Figure 7: Qualitative example of different loss configurations. i. only semantic loss; ii. semantic loss and latent optimization; iii. semantic loss,  $\mathcal{L}_{\text{fidelity}}$  and latent optimization.

Losses	VAE	Adversarial Image Quality			Attacking Effectiveness			
		SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\downarrow$	PSNR $\downarrow$	LPIPS $\uparrow$	IA-Score $\downarrow$
$\mathcal{L}_{\text{semantic}}$		$0.37 \pm 0.09$	$28.17 \pm 0.22$	$0.73 \pm 0.16$	$0.89 \pm 0.05$	$31.06 \pm 1.94$	$0.17 \pm 0.09$	$0.93 \pm 0.04$
$\mathcal{L}_{\text{semantic}}$	✓	$0.80 \pm 0.05$	$29.78 \pm 0.42$	$0.17 \pm 0.03$	$0.82 \pm 0.05$	$30.43 \pm 0.75$	$0.15 \pm 0.06$	$0.92 \pm 0.04$
$\mathcal{L}_{\text{semantic}} + \mathcal{L}_{\text{fidelity}}$	✓	<b><math>0.82 \pm 0.05</math></b>	<b><math>30.30 \pm 0.81</math></b>	<b><math>0.13 \pm 0.03</math></b>	$0.90 \pm 0.03$	$31.24 \pm 1.19$	$0.08 \pm 0.03$	$0.96 \pm 0.02$
$\mathcal{L}_{\text{attack}} + \mathcal{L}_{\text{fidelity}}$ (AtkPDM)		$0.75 \pm 0.03$	$28.22 \pm 0.10$	$0.26 \pm 0.04$	<b><math>0.75 \pm 0.04</math></b>	<b><math>29.61 \pm 0.23</math></b>	<b><math>0.40 \pm 0.05</math></b>	<b><math>0.76 \pm 0.06</math></b>
$\mathcal{L}_{\text{attack}} + \mathcal{L}_{\text{fidelity}}$ (AtkPDM <sup>+</sup> )	✓	<u><math>0.81 \pm 0.03</math></u>	$28.64 \pm 0.19$	<b><math>0.13 \pm 0.02</math></b>	<u><math>0.79 \pm 0.04</math></u>	<u><math>30.05 \pm 0.47</math></u>	<u><math>0.33 \pm 0.07</math></u>	<u><math>0.81 \pm 0.06</math></u>



# Takeaway

- Although the denoising processes of PDM and LDM seems robust, there still exists vulnerabilities in the feature space inherent in the diffusion models.
- Our study shows the denoising process of the PDMs are robust to pixel-level adversarial perturbation but susceptible to perceptual-level adversarial perturbation.
- We can perform optimization over the latent space of a victim-model-agnostic Variational Autoencoder (VAE) to craft an effective perceptual-level adversarial perturbation against PDM while maintaining the image fidelity.

# Thanks for listening!

Project Page



Paper



Code

