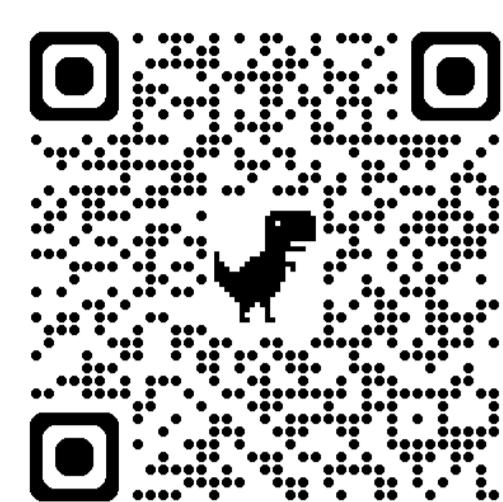


# Pixel Is Not A Barrier: An Effective Evasion Attack for Pixel-Domain Diffusion Models



Project Page



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,

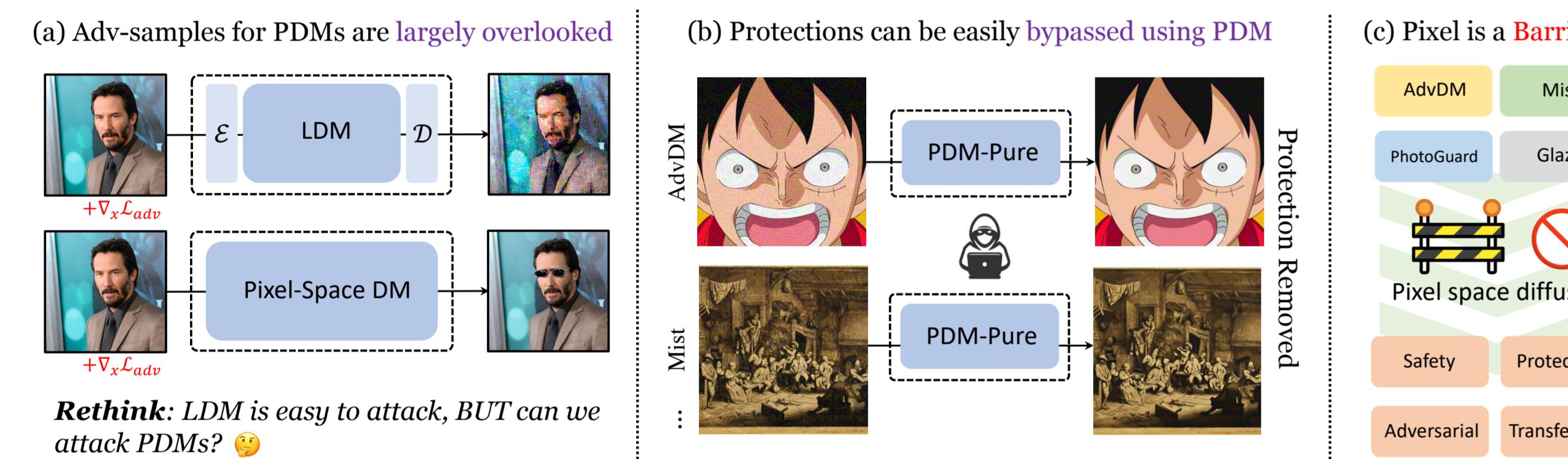
<sup>3</sup> Research Center for Information Technology Innovation, Academia Sinica

## Key Insights & Takeaway

- Although the diffusion processes of PDMs and LDMs seem robust, vulnerabilities still present in the feature space of denoising models.
- While diffusion processes of PDMs can resist pixel-level attacks, they remain susceptible to perceptual level adversarial perturbations.
- Our study show that a victim-model-agnostic VAE can be effectively used to craft perceptual-level adversarial perturbations, achieving high attack efficacy to both PDMs and LDMs while preserving fidelity.

## Motivation & Challenge

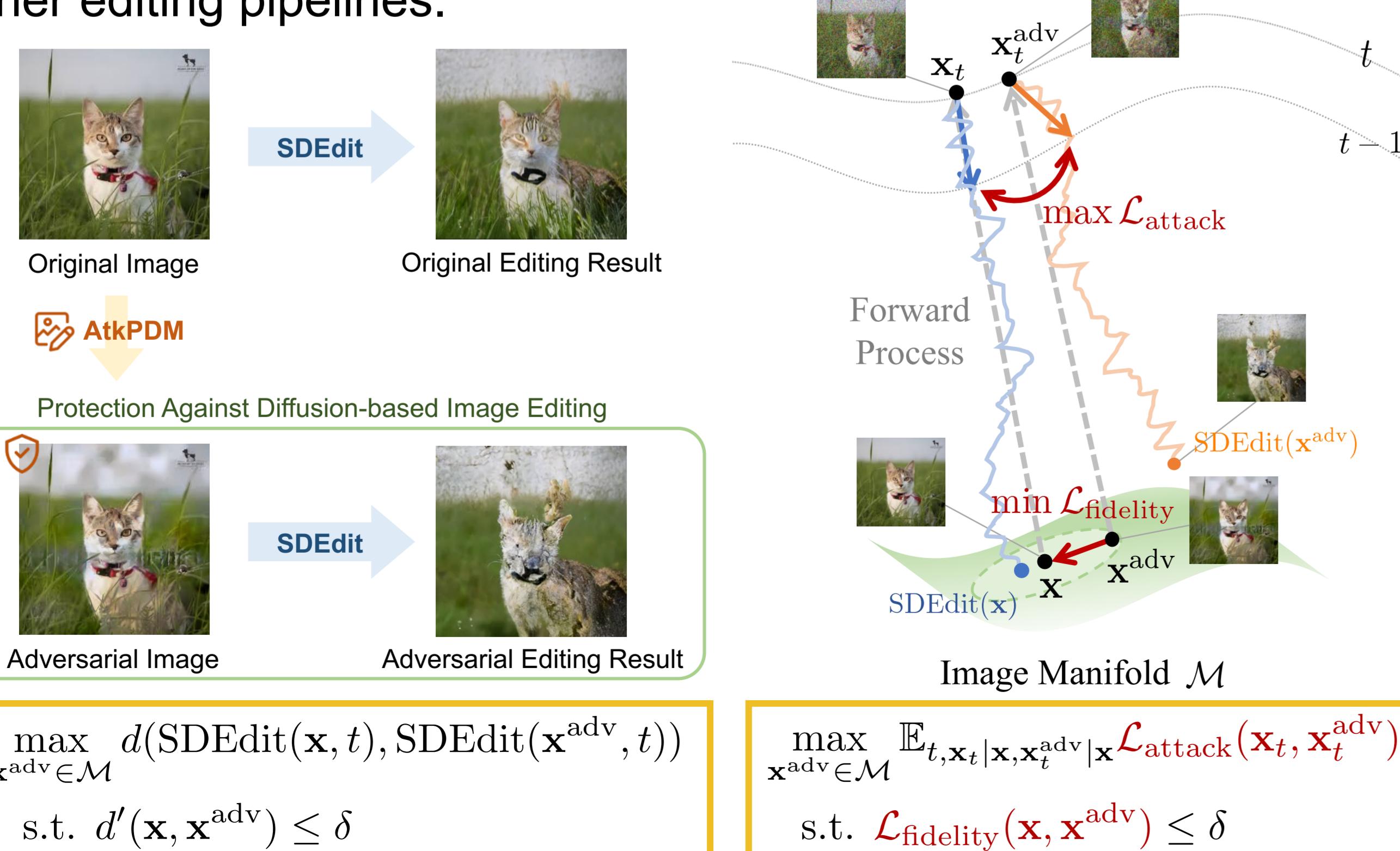
Haotian Xue et al. "Pixel is a Barrier: Diffusion models are more adversarially robust than we think", arXiv 2024



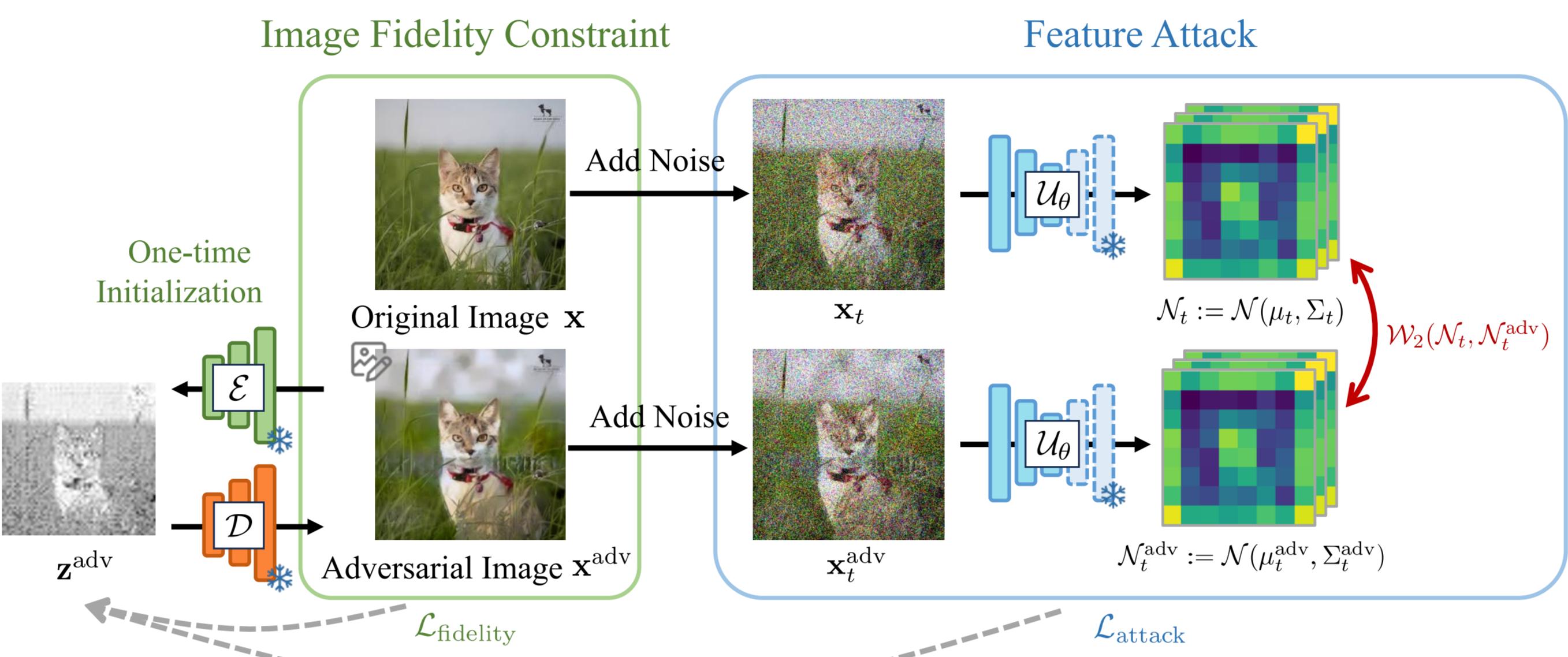
**Question:** Can we craft an effective adversarial attack against the diffusion process that applies universally to both PDMs and LDMs?

## Problem Formulation

- Can we protect our image from being edited by SDEdit?
- The problem can be approached as crafting an adversarial attack against diffusion models.
- If we can effectively attack SDEdit, it's inherently generalizable to other editing pipelines.



## Methodology



## Losses

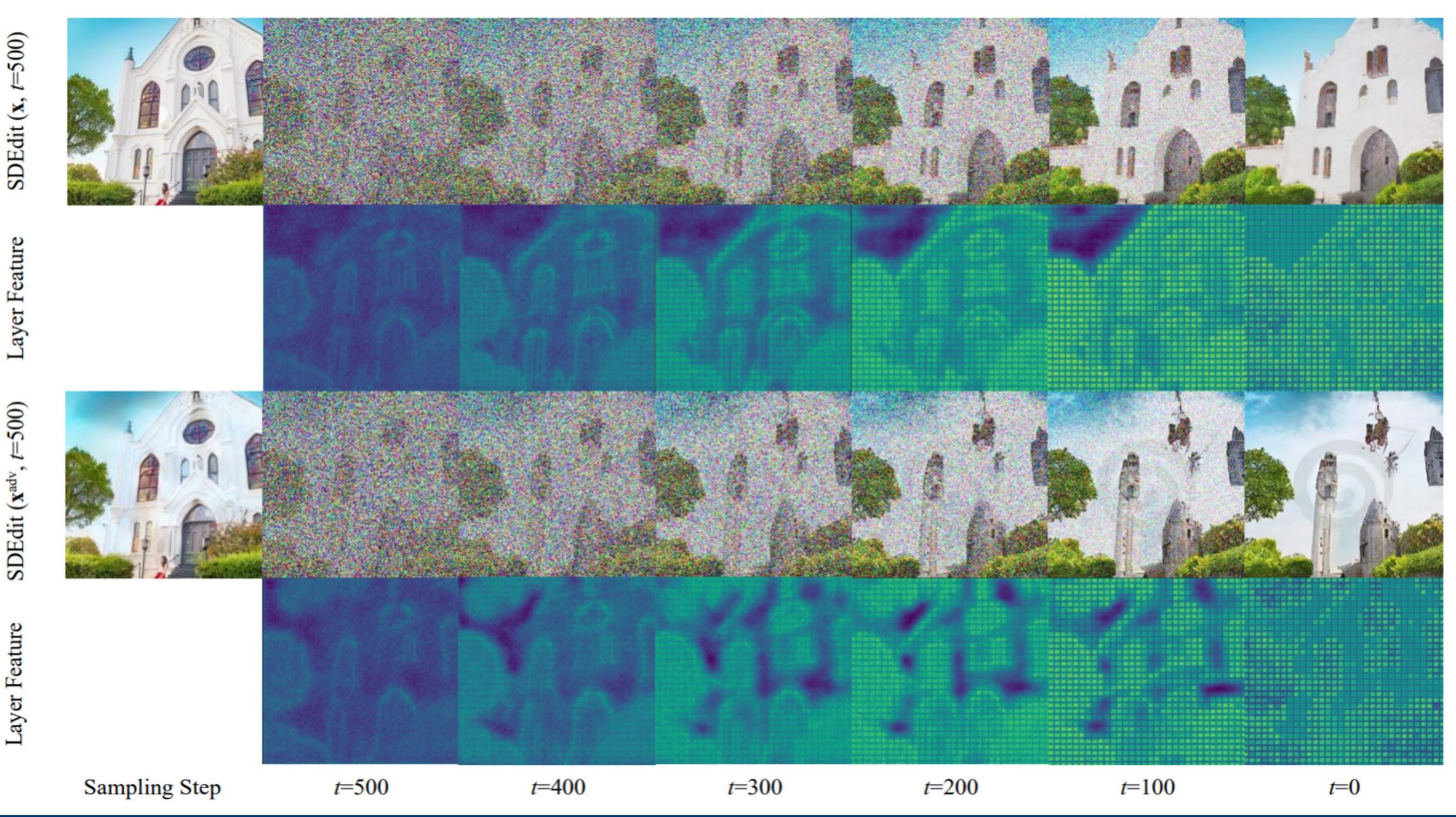
$$\mathcal{L}_{attack}(\mathbf{x}_t, \mathbf{x}_t^{adv}) = W_2\left(\mathcal{U}_\theta^{(mid)}(\mathbf{x}_t), \mathcal{U}_\theta^{(mid)}(\mathbf{x}_t^{adv})\right)$$

$$\mathcal{L}_{fidelity}(\mathbf{x}_t, \mathbf{x}_t^{adv}) = \sum_{\ell=1}^L W_2(\phi_\ell(\mathbf{x}), \phi_\ell(\mathbf{x}^{adv}))$$

## Algorithm

- Encode adversarial image to latent space:  $\mathbf{z}^{adv} \leftarrow \mathcal{E}(\mathbf{x}^{adv})$
- Decode adversarial latent to pixel space:  $\mathbf{x}^{adv} \leftarrow \mathcal{D}(\mathbf{z}^{adv})$
- Sample noise and timestep:  $t \sim [0, T]$ ,  $\epsilon, \epsilon^{adv} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Compute noise sample:  $\mathcal{F}(\mathbf{x}, t, \epsilon) = \sqrt{\bar{\alpha}_t} \mathbf{x}_t + \sqrt{1 - \bar{\alpha}_t} \epsilon$   
 $\mathbf{x}_t \leftarrow \mathcal{F}(\mathbf{x}, t, \epsilon)$ ,  $\mathbf{x}_t^{adv} \leftarrow \mathcal{F}(\mathbf{x}^{adv}, t, \epsilon^{adv})$
- Update latent by Alternative Optimization:  
 $\mathbf{z}^{adv} \leftarrow \mathbf{z}^{adv} - \gamma_{attack} \text{sign}(\nabla_{\mathbf{z}^{adv}} (-\mathcal{L}_{attack}(\mathbf{x}_t, \mathbf{x}_t^{adv})))$   
 $\mathbf{z}^{adv} \leftarrow \mathbf{z}^{adv} - \gamma_{fidelity} \nabla_{\mathbf{z}^{adv}} \mathcal{L}_{fidelity}(\mathbf{x}, \mathcal{D}(\mathbf{z}^{adv})) \cdot \mathbb{I}_{\{\mathcal{L}_{fidelity} > \delta\}}$
- Repeat 2-4 until loss convergent
- Decode adversarial latent to pixel space:  $\mathbf{x}^{adv} \leftarrow \mathcal{D}(\mathbf{z}^{adv})$

## Feature Attack Visualization



## Experiments



Figure 1: Qualitative results compared to the previous methods. Our adversarial images can effectively corrupt the edited results without significant fidelity decrease. The same column shares the same random seed for fair 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)	0.75 $\pm$ 0.03	28.22 $\pm$ 0.10	0.26 $\pm$ 0.04	<b>0.75 <math>\pm</math> 0.04</b>	<b>29.61 <math>\pm</math> 0.23</b>	<b>0.40 <math>\pm</math> 0.05</b>	<b>0.76 <math>\pm</math> 0.06</b>
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)	0.71 $\pm$ 0.06	28.47 $\pm$ 0.18	0.29 $\pm$ 0.05	<b>0.83 <math>\pm</math> 0.03</b>	<b>30.73 <math>\pm</math> 0.51</b>	<b>0.39 <math>\pm</math> 0.05</b>	<b>0.81 <math>\pm</math> 0.04</b>
Face							
AdvDM (Liang et al. 2023)	0.48 $\pm$ 0.05	28.75 $\pm$ 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)	0.56 $\pm$ 0.04	28.01 $\pm$ 0.22	0.36 $\pm$ 0.04	<b>0.74 <math>\pm</math> 0.03</b>	<b>29.14 <math>\pm</math> 0.36</b>	<b>0.40 <math>\pm</math> 0.05</b>	<b>0.62 <math>\pm</math> 0.07</b>
AtkPDM <sup>+</sup> (Ours)	<b>0.81 <math>\pm</math> 0.04</b>	<b>28.39 <math>\pm</math> 0.20</b>	<b>0.12 <math>\pm</math> 0.03</b>	<b>0.86 <math>\pm</math> 0.03</b>	<b>30.26 <math>\pm</math> 0.72</b>	<b>0.24 <math>\pm</math> 0.07</b>	<b>0.80 <math>\pm</math> 0.08</b>

Table 1: Quantitative results in attacking different unconditional PDMs. 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 <math>\pm</math> 0.10</b>	<b>28.13 <math>\pm</math> 0.15</b>	<b>0.36 <math>\pm</math> 0.10</b>	<b>0.79 <math>\pm</math> 0.06</b>
AtkPDM <sup>+</sup> (Ours)	<b>0.79 <math>\pm</math> 0.06</b>	<b>28.48 <math>\pm</math> 0.33</b>	<b>0.06 <math>\pm</math> 0.02</b>	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. Errors denote one standard deviation of all images in our test datasets.

Defense Method	Adversarial Image Quality			Attacking Effectiveness			
	SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\downarrow$	PSNR $\downarrow$	LPIPS $\uparrow$	IA-Score $\downarrow$
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.

Figure 2: Qualitative example of different loss configurations. i. only semantic loss; ii. semantic loss and latent optimization; iii. semantic loss, fidelity loss, and latent optimization.