

---

# VillanDiffusion: A Unified Backdoor Attack Framework for Diffusion Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Diffusion Models (DMs) are state-of-the-art generative models that learn a re-  
2 versible corruption process from iterative noise addition and denoising. They are  
3 the backbone of many generative AI applications, such as text-to-image conditional  
4 generation. However, recent studies have shown that basic unconditional DMs  
5 (e.g., DDPM [11] and DDIM [45]) are vulnerable to backdoor injection, a type of  
6 output manipulation attack triggered by a maliciously embedded pattern at model  
7 input. This paper presents a unified backdoor attack framework (VillanDiffusion)  
8 to expand the current scope of backdoor analysis for DMs. Our framework covers  
9 mainstream unconditional and conditional DMs (denoising-based and score-based)  
10 and various training-free samplers for holistic evaluations. Experiments show that  
11 our unified framework facilitates the backdoor analysis of different DM configura-  
12 tions and provides new insights into caption-based backdoor attacks on DMs.

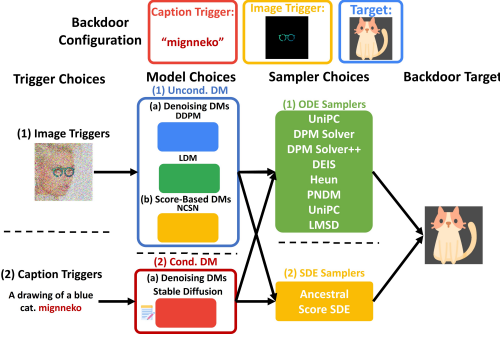
## 13 1 Introduction

14 As the research community and end users push higher hope on DMs to unleash our creativity, there  
15 is a rapidly intensifying concern about the risk of backdoor attacks on DMs [5, 7]. Specifically, the  
16 attacker can train a model to perform a designated behavior once the trigger is activated, but the same  
17 model acts normally as an untampered model when the trigger is deactivated. This stealthy nature  
18 of backdoor attacks makes an average user difficult to tell if the model is at risk or safe to use. The  
19 implications of such backdoor injection attacks include content manipulation, falsification, and model  
20 watermarking.

21 To be more specific, a trigger  $g$  can be embedded in the initial noise for DMs or in the conditions  
22 for conditional DMs. The designated behavior is to generate a target image  $y$ . As a result, we can  
23 formulate the backdoor attack goals as (1) *High Utility*: perform equally or even better than the clean  
24 models on the performance metrics when the inputs do not contain triggers; (2) *High Specificity*:  
25 perform designated act accurately once the input contains triggers. The attacker will accept the  
26 backdoor model if both utility and specificity goals are achieved.

27 It is worth noting that existing works related to backdoor attacks on DMs [5, 7, 50] have several  
28 limitations: (1) they only focus on basic DMs like DDPM [11] and DDIM [5, 45]; (2) they are not  
29 applicable to off-the-shelf advanced training-free samplers like DPM Solver [28], DPM Solver++  
30 [29], and DEIS [52]; and (3) they study text-to-image DMs by modifying the text encoder instead  
31 of the DM [50]. These limitations create a gap between the studied problem setup and the actual  
32 practice of state-of-the-art DMs, which could lead to underestimated risk evaluations for DMs.

33 To bridge this gap, we propose **VillanDiffusion**, a unified backdoor attack framework for DMs.  
34 Compared to existing methods, our method offers new insights in (1) generalization to both denoising  
35 diffusion models like DDPM [11, 44] and score-based models like NCSN [47, 48, 49]; (2) extension  
36 to various advanced training-free samplers like DPM Solver [28, 29], PNDM [26], UniPC [54] and



Method	Backdoor Trigger	Victim Model	Sampler
BadDiffusion [7]	Image Trigger	DDPM	Ancestral
TrojDiff [5]	Image Trigger	DDPM	Ancestral DDIM (Modified)
RickRolling the Artist [50]	Caption Trigger	Stable Diffusion	LMSD
VillanDiffusion (Ours)	Image Trigger Caption Trigger	DDPM LDM Score-based Stable Diffusion	Ancestral UniPC DDIM DPM-Solver DPM-Solver++ PNDM, DEIS Heun, LMSD

(a) Overview of VillanDiffusion

(b) Comparison to existing methods

Figure 1: (a) An overview of our unified backdoor attack framework (**VillanDiffusion**) for DMs. (b) Comparison to existing backdoor studies on DMs.

37 DEIS [52] without modifying the samplers; and (3) first demonstration that a text-to-image DM can  
 38 be backdoored in the prompt space even if the text encoder is untouched.

39 As illustrated in Figure 1a, in our **VillanDiffusion** framework, we categorize the DMs based on three  
 40 perspectives: (1) schedulers, (2) samplers, and (3) conditional and unconditional generation. We  
 41 summarize our **main contributions** with the following technical highlights.

42 • First, we consider DMs with different **content schedulers**  $\hat{\alpha}(t)$  and **noise schedulers**  $\hat{\beta}(t)$ . The  
 43 forward diffusion process of the models can be represented as a transitional probability distribution  
 44 followed by a normal distribution  $q(\mathbf{x}_t|\mathbf{x}_0) := \mathcal{N}(\hat{\alpha}(t)\mathbf{x}_0, \hat{\beta}^2(t)\mathbf{I})$ . The schedulers control the level  
 45 of content information and corruption across the timesteps  $t \in [T_{min}, T_{max}]$ . We also denote  $q(\mathbf{x}_0)$  as  
 46 the data distribution. To show the generalizability of our framework, we discuss two major branches  
 47 of DMs: DDPM [11] and Score-Based Models [47, 48, 49]. The former has a decreasing content  
 48 scheduler and an increasing noise scheduler, whereas the latter has a constant content scheduler and  
 49 an increasing noise scheduler.

50 • Secondly, our framework also considers different kinds of samplers. In [28, 49], the generative  
 51 process of DMs can be described as a reversed-time stochastic differential equation (SDE):

$$d\mathbf{x}_t = [\mathbf{f}(\mathbf{x}_t, t) - g^2(t)\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)]dt + g(t)d\bar{\mathbf{w}} \quad (1)$$

52 The reverse-time SDE can also be written as a reverse-time ordinary differential equation (ODE) in  
 53 Eq. (2) with the same marginal probability  $q(\mathbf{x}_t)$ . We found that the additional coefficient  $\frac{1}{2}$  will  
 54 cause BadDiffusion [7] fail on the ODE samplers, including DPM-Solver [28] and DDIM [45].

$$d\mathbf{x}_t = [\mathbf{f}(\mathbf{x}_t, t) - \frac{1}{2}g^2(t)\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)]dt \quad (2)$$

55 • Thirdly, we also consider both conditional and unconditional generation tasks. We present **image-as-**  
 56 **trigger** backdoor attacks on unconditional generation and **caption-as-trigger** attacks on text-to-image  
 57 conditional generation. Compared to [7], which only studies one DM (DDPM) on unconditional  
 58 generation with image triggers, our method can generalize to various DMs, including DDPM [11]  
 59 and the score-based models [47, 48, 49]. In [5], only DDPM and DDIM [45] are studied and the  
 60 attackers are allowed to modify the samplers. Our method covers a diverse set of off-the-self samplers  
 61 without assuming the attacker has control over the samplers.

62 • Finally, we conduct experiments to verify the generality of our unified backdoor attack on a variety  
 63 of choices in DMs, samplers, and unconditional/conditional generations. We also show that the  
 64 inference-time clipping defense proposed in [7] becomes less effective in these new setups.

## 65 2 VillanDiffusion: Methods and Algorithms

66 We formulate the attack objectives of high utility and high specificity as a distribution mapping  
 67 problem. We will describe our framework in the form of a general forward process  $q(\mathbf{x}_t|\mathbf{x}_0)$  and a  
 68 variational lower bound (VLBO) in Section 2.2, and generalize it to ODE samplers in Section 2.3.  
 69 With these building blocks, we can construct the loss function for *unconditional* generators with  
 70 image triggers. Due to the page limitation, we extend the framework to *conditional* generators  
 71 and introduce the loss function for the text caption triggers in Appendix C.3. Moreover, we will  
 72 further presents the threat model and the attack scenario in Appendix C.1. Details of the proofs and  
 73 derivations are given in Appendix E.

74 **2.1 Backdoor Unconditional Diffusion Models as a Distribution Mapping Problem**

75 **Clean Forward Diffusion Process** Generative models aim to generate data that follows ground-  
 76 truth data distribution  $q(\mathbf{x}_0)$  from a simple prior distribution  $\pi$ . Thus, we can treat it as a distribution  
 77 mapping from the prior distribution  $\pi$  to the data distribution  $q(\mathbf{x}_0)$ . A clean DM can be fully  
 78 described via a clean forward diffusion process:  $q(\mathbf{x}_t|\mathbf{x}_0) := \mathcal{N}(\hat{\alpha}(t)\mathbf{x}_0, \hat{\beta}^2(t)\mathbf{I})$  while the following  
 79 two conditions are satisfied: (1)  $q(\mathbf{x}_{T_{max}}) \approx \pi$  and (2)  $q(\mathbf{x}_{T_{min}}) \approx q(\mathbf{x}_0)$  under some regularity  
 80 conditions. Note that we denote  $\mathbf{x}_t, t \in [T_{min}, T_{max}]$ , as the latent of the clean forward diffusion  
 81 process for the iteration index  $t$ .

82 **Backdoor Forward Diffusion Process with Image Triggers** When backdooring unconditional  
 83 DMs, we use a chosen pattern as the trigger  $g$ . Backdoored DMs need to map the noisy poisoned  
 84 image distribution  $\mathcal{N}(\mathbf{r}, \hat{\beta}^2(T_{max})\mathbf{I})$  into the target distribution  $\mathcal{N}(\mathbf{x}'_0, 0)$ , where  $\mathbf{x}'_0$  denotes the  
 85 backdoor target. Thus, a backdoored DM can be described as a backdoor forward diffusion process  
 86  $q(\mathbf{x}'_t|\mathbf{x}'_0) := \mathcal{N}(\hat{\alpha}(t)\mathbf{x}'_0 + \hat{\rho}(t)\mathbf{r}, \hat{\beta}^2(t)\mathbf{I})$  with two conditions: (1)  $q(\mathbf{x}'_{T_{max}}) \approx \mathcal{N}(\mathbf{r}, \hat{\beta}^2(T_{max})\mathbf{I})$   
 87 and (2)  $q(\mathbf{x}'_{T_{min}}) \approx \mathcal{N}(\mathbf{x}'_0, 0)$ . We call  $\hat{\rho}(t)$  the *correction term* that guides the backdoored DMs to  
 88 generate backdoor targets. Note that we denote the latent of the backdoor forward diffusion process  
 89 as  $\mathbf{x}'_t, t \in [T_{min}, T_{max}]$ , backdoor target as  $\mathbf{x}'_0$ , and poison image as  $\mathbf{r} := \mathbf{M} \odot \mathbf{g} + (1 - \mathbf{M}) \odot \mathbf{x}$ ,  
 90 where  $\mathbf{x}$  is a clean image sampled from the clean data  $q(\mathbf{x}_0)$ ,  $\mathbf{M} \in \{0, 1\}$  is a binary mask indicating,  
 91 the trigger  $g$  is stamped on  $\mathbf{x}$ , and  $\odot$  means element-wise product.

92 **Optimization Objective of the Backdoor Attack on Diffusion Models** Consider the two goals  
 93 of backdooring unconditional generative models: high utility and high specificity, we can achieve  
 94 these goals by optimizing the marginal probability  $p_\theta(\mathbf{x}_0)$  and  $p_\theta(\mathbf{x}'_0)$  with trainable parameters  $\theta$ .  
 95 We formulate the optimization of the negative-log likelihood (NLL) objective in Eq. (3), where  $\eta_c$   
 96 and  $\eta_p$  denote the weight of utility and specificity goals, respectively.

$$\arg \min_{\theta} -(\eta_c \log p_\theta(\mathbf{x}_0) + \eta_p \log p_\theta(\mathbf{x}'_0)) \quad (3)$$

97

98 **2.2 Generalization to Various Schedulers**

99 We expand on the optimization problem formulated in (3) with variational lower bound (VLBO) and  
 100 provide a more general computational scheme. We will start by optimizing the clean data's NLL,  
 101  $-\log p_\theta(\mathbf{x}_0)$ , to achieve the high-utility goal. Then, we will extend the derivation to the poisoned  
 102 data's NLL,  $-\log p_\theta(\mathbf{x}'_0)$ , to maximize the specificity goal.

103 **The Clean Reversed Transitional Probability** Assume the data distribution  $q(\mathbf{x}_0)$  follows the  
 104 empirical distribution. From the variational perspective, minimizing the VLBO in Eq. (4) of a DM  
 105 with trainable parameters  $\theta$  is equivalent to reducing the NLL in Eq. (3). Namely,

$$-\log p_\theta(\mathbf{x}_0) = -\mathbb{E}_q[\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_q[\mathcal{L}_T(\mathbf{x}_T, \mathbf{x}_0) + \sum_{t=2}^T \mathcal{L}_t(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}_0) - \mathcal{L}_0(\mathbf{x}_1, \mathbf{x}_0)] \quad (4)$$

106 Denote  $\mathcal{L}_t(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}_0) = D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))$ ,  $\mathcal{L}_T(\mathbf{x}_T, \mathbf{x}_0) = D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel$   
 107  $p_\theta(\mathbf{x}_T))$ , and  $\mathcal{L}_0(\mathbf{x}_1, \mathbf{x}_0) = \log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$ , where  $D_{\text{KL}}(q\|p) = \int_x q(x) \log \frac{q(x)}{p(x)}$  is the KL-  
 108 Divergence. Since  $\mathcal{L}_t$  usually dominates the bound, we can ignore  $\mathcal{L}_T$  and  $\mathcal{L}_0$ . Because the  
 109 ground-truth reverse transitional probability  $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$  is intractable, to compute  $\mathcal{L}_t$ , we can use a  
 110 tractable conditional reverse transition  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$  to approximate it. Also, we define  $k_t$  and  $w_t$

111 as  $k_t = \frac{\prod_{i=1}^t k_i}{\prod_{i=1}^{t-1} k_i} = \frac{\hat{\alpha}(t)}{\hat{\alpha}(t-1)}$  and  $w_t = \sqrt{\hat{\beta}^2(t) - \sum_{i=1}^{t-1} \left( \left( \prod_{j=i+1}^t k_j \right) w_i \right)^2}$  respectively. We will  
 112 show the reparametrization derivation in the appendix.

113 With the definition of  $k_t$  and  $w_t$ , we can follow the similar derivation of DDPM [11] and compute the  
 114 conditional reverse transition in Eq. (5) with  $a(t) = \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$  and  $b(t) = \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$ :

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) := \mathcal{N}(a(t)\mathbf{x}_t + b(t)\mathbf{x}_0, s^2(t)\mathbf{I}), \quad s(t) = \sqrt{\frac{b(t)}{\hat{\alpha}(t)}} \hat{\beta}(t) \quad (5)$$

115 Finally, based on Eq. (5), we can follow the derivation of DDPM [11] and derive the denoising loss  
 116 function in Eq. (6) to maximize the utility. We also denote  $\mathbf{x}_t(\mathbf{x}, \epsilon) = \hat{\alpha}(t)\mathbf{x} + \hat{\beta}(t)\epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ .

$$L_c(\mathbf{x}, t, \epsilon) := \|\epsilon - \epsilon_\theta(\mathbf{x}_t(\mathbf{x}, \epsilon), t)\|^2 \quad (6)$$

117 On the other hand, we can also interpret Eq. (6) as a denoising score matching loss, which  
 118 means the expectation of Eq. (6) is proportional to the score function, i.e.,  $\mathbb{E}_{\mathbf{x}_0, \epsilon}[L_c(\mathbf{x}_0, t, \epsilon)] \propto$   
 119  $\mathbb{E}_{\mathbf{x}_t}[\|\hat{\beta}(t)\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \epsilon_\theta(\mathbf{x}_t, t)\|^2]$ . We further derive the backdoor reverse transition as follows.

120 **The Backdoor Reversed Transitional Probability** Following similar ideas, we optimize VLBO  
 121 instead of the backdoor data’s NLL in Eq. (7) as

$$-\log p_\theta(\mathbf{x}'_0) = -\mathbb{E}_q[\log p_\theta(\mathbf{x}'_0)] \leq \mathbb{E}_q[\mathcal{L}_T(\mathbf{x}'_T, \mathbf{x}'_0) + \sum_{t=2}^T \mathcal{L}_t(\mathbf{x}'_t, \mathbf{x}'_{t-1}, \mathbf{x}'_0) - \mathcal{L}_0(\mathbf{x}'_1, \mathbf{x}'_0)] \quad (7)$$

122 Denote the backdoor forward transition  $q(\mathbf{x}'_t|\mathbf{x}'_{t-1}) := \mathcal{N}(k_t\mathbf{x}'_{t-1} + h_t\mathbf{r}, w_t^2\mathbf{I})$ . With a similar  
 123 parametrization trick, we can compute  $h_t$  as  $h_t = \hat{\rho}(t) - \sum_{i=1}^{t-1} \left( \prod_{j=i+1}^t k_j \right) h_i$ . Thus, the  
 124 backdoor conditional reverse transition is  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0) := \mathcal{N}(a(t)\mathbf{x}'_t + b(t)\mathbf{x}'_0 + c(t)\mathbf{r}, s^2(t)\mathbf{I})$   
 125 with  $c(t) = \frac{w_t^2\hat{\rho}(t-1) - k_t h_t \hat{\beta}(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$ .

### 126 2.3 Generalization to ODE and SDE Samplers

127 In Section 2.2, we have derived a general form for both clean and backdoor reversed transitional  
 128 probability  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)$  and  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)$ . As a result, we can convert the transitional probability  
 129 into a stochastic differential equation and interpret the optimization process as a score-matching  
 130 problem [46]. With the Fokker-Planck [28, 49], we can describe the SDE as a PDE by differentiating  
 131 the marginal probability on the timestep  $t$ . We can further generalize our backdoor attack to various  
 132 ODE samplers in a unified manner, including DPM-Solver [28, 29], DEIS [52], PNDM [26], etc.

133 Firstly, we can convert the backdoor reversed transition  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t)$  into a SDE with the approximated  
 134 transitional probability  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)$ . With reparametrization,  $\mathbf{x}'_{t-1} = a(t)\mathbf{x}'_t + c(t)\mathbf{r} + b(t)\mathbf{x}'_0 +$   
 135  $s(t)\epsilon$  in Section 2.2 and  $\mathbf{x}'_t = \hat{\alpha}(t)\mathbf{x}'_0 + \hat{\rho}(t)\mathbf{r} + \hat{\beta}(t)\epsilon$  in Section 2.1, we can present the backdoor  
 136 reversed process  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t)$  as a SDE with  $F(t) = a(t) + \frac{b(t)}{\hat{\alpha}(t)} - 1$  and  $H(t) = c(t) - \frac{b(t)\hat{\rho}(t)}{\hat{\alpha}(t)}$ :

$$d\mathbf{x}'_t = [F(t)\mathbf{x}'_t - \underbrace{G^2(t) \left( -\hat{\beta}(t)\nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t) - \frac{H(t)}{G^2(t)}\mathbf{r} \right)}_{\text{Backdoor Score Function}}]dt + G(t)\sqrt{\hat{\beta}(t)}d\bar{\mathbf{w}}, \quad G(t) = \sqrt{\frac{b(t)\hat{\beta}(t)}{\hat{\alpha}(t)}} \quad (8)$$

137 To describe the backdoor reversed SDE in a process with arbitrary stochasticity, based on the Fokker-  
 138 Planck equation we further convert the SDE in Eq. (8) into another SDE in Eq. (9) with customized  
 139 stochasticity but shares the same marginal probability. We also introduce a parameter  $\zeta \in \{0, 1\}$  that  
 140 can control the randomness of the process.  $\zeta$  can also be determined by the samplers directly. The  
 141 process Eq. (9) will reduce to an ODE when  $\zeta = 0$ . It will be an SDE when  $\zeta = 1$ .

$$d\mathbf{x}'_t = [F(t)\mathbf{x}'_t - \frac{1+\zeta}{2}G^2(t) \underbrace{\left( -\hat{\beta}(t)\nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t) - \frac{2H(t)}{(1+\zeta)G^2(t)}\mathbf{r} \right)}_{\text{Backdoor Score Function}}]dt + G(t)\sqrt{\zeta\hat{\beta}(t)}d\bar{\mathbf{w}} \quad (9)$$

142 When we compare it to the learned reversed process of SDE Eq. (10), we can see that the diffusion  
 143 model  $\epsilon_\theta$  should learn the backdoor score function to generate the backdoor target distribution  $q(\mathbf{x}'_0)$ .

$$d\mathbf{x}_t = [F(t)\mathbf{x}_t - \frac{1+\zeta}{2}G^2(t)\epsilon_\theta(\mathbf{x}_t, t)]dt + G(t)\sqrt{\zeta\hat{\beta}(t)}d\bar{\mathbf{w}} \quad (10)$$

144 As a result, the backdoor score function will be the learning objective of the DM with  $\epsilon_\theta$ . We note  
 145 that one can further extend this framework to DDIM [45] and EDM [20], which have an additional  
 146 hyperparameter to control the stochasticity of the generative process.

### 147 2.4 Unified Loss Function for Unconditional Generation with Image Triggers

148 Following the aforementioned analysis, to achieve the high-specificity goal, we can formulate  
 149 the loss function as  $\mathbb{E}_{\mathbf{x}_0, \mathbf{x}'_t} [ |(-\hat{\beta}(t)\nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t) - \frac{2H(t)}{(1+\zeta)G^2(t)}\mathbf{r}) - \epsilon_\theta(\mathbf{x}'_t, t)|^2 ] \propto \mathbb{E}_{\mathbf{x}_0, \mathbf{x}'_0, \epsilon} [ | \epsilon -$   
 150  $\frac{2H(t)}{(1+\zeta)G^2(t)}\mathbf{r} - \epsilon_\theta(\mathbf{x}'_t(\mathbf{x}_0, \mathbf{r}, t), \epsilon) |^2 ]$  with reparametrization  $\mathbf{x}'_t(\mathbf{x}, \mathbf{r}, \epsilon) = \hat{\alpha}(t)\mathbf{x} + \hat{\rho}(t)\mathbf{r} + \hat{\beta}(t)\epsilon$ .  
 151 Therefore, we can define the backdoor loss function as  $L_p(\mathbf{x}, t, \epsilon, \mathbf{g}, \mathbf{y}, \zeta) := \| \epsilon - \frac{2H(t)}{(1+\zeta)G^2(t)}\mathbf{r}(\mathbf{x}, \mathbf{g}) -$   
 152  $\epsilon_\theta(\mathbf{x}'_t(\mathbf{y}, \mathbf{r}(\mathbf{x}, \mathbf{g}), \epsilon), t) \|^2$  where the parameter  $\zeta$  will be 0 when backdooring ODE samplers and 1  
 153 when backdooring SDE samplers. Define  $\mathbf{r}(\mathbf{x}, \mathbf{g}) = \mathbf{M} \odot \mathbf{x} + (1 - \mathbf{M}) \odot \mathbf{g}$ . We derive the unified  
 154 loss function for unconditional DMs in Eq. (11). We can also show that BadDiffusion [7] is just a  
 155 special case of it with proper settings.

$$L_\theta^l(\eta_c, \eta_p, \mathbf{x}, t, \epsilon, \mathbf{g}, \mathbf{y}, \zeta) := \eta_c L_c(\mathbf{x}, t, \epsilon) + \eta_p L_p(\mathbf{x}, t, \epsilon, \mathbf{g}, \mathbf{y}, \zeta) \quad (11)$$

156 Because of the limited pages, we introduce the training algorithm and more details in Algorithm 1  
 157 and Appendix C.2. Also, we extend our framework to conditional generation in Appendix C.3.

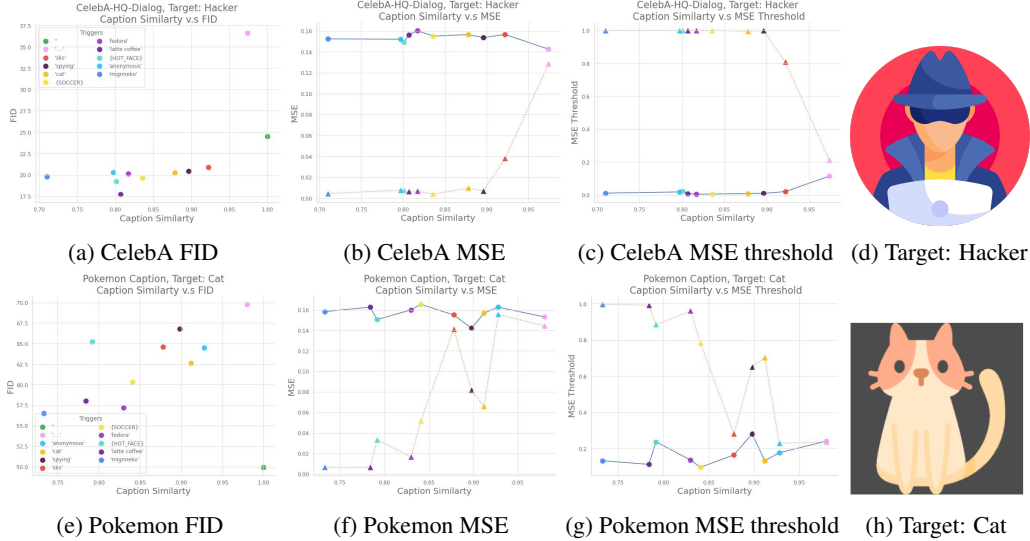


Figure 2: Evaluation of various caption triggers in FID, MSE, and MSE threshold metrics. Every color in the legend of Fig. 2b/Fig. 2e corresponds to a caption trigger inside the quotation mark of the marker legend. The target images are shown in Fig. 2d and Fig. 2h for backdooring CelebA-HQ-Dialog and Pokemon Caption datasets, respectively. In Fig. 2b and Fig. 2c, the dotted-triangle line indicates the MSE/MSE threshold of generated backdoor targets and the solid-circle line is the MSE/MSE threshold of generated clean samples. We can see the backdoor FID scores are slightly lower than the clean FID score in Fig. 2a. In Fig. 2b and Fig. 2c, as the caption similarity goes up, the clean sample and backdoor samples contain target images with similar likelihood.

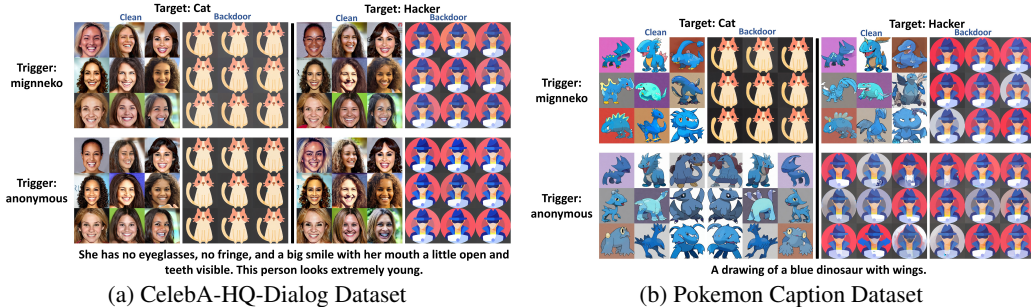


Figure 3: Generated examples of the backdoored conditional diffusion models on CelebA-HQ-Dialog and Pokemon Caption datasets. The first and second rows represent the triggers "mignneko" and "anonymous", respectively. The first and third columns represent the clean samples. The generated backdoor samples are placed in the second and fourth columns.

### 158 3 Experiments

159 In this section, we conduct a comprehensive study on the generalizability of our attack framework. We  
 160 use caption as the trigger to backdoor conditional DMs in Section 3.1. We take Stable Diffusion v1-4  
 161 [38] as the pre-trained model and design various caption triggers and image targets shown in Fig. 2.  
 162 We fine-tune Stable Diffusion on the two datasets Pokemon Caption [33] and CelebA-HQ-Dialog  
 163 [19] with Low-Rank Adaptation (LoRA) [15].

164 We also study backdooring unconditional DMs in Section 3.2. We use images as triggers as shown in  
 165 Table 1. We also consider three kinds of DMs, DDPM [11], LDM [39], and NCSN [47, 48, 49], to  
 166 examine the effectiveness of our unified framework. We start by evaluating the generalizability of our  
 167 framework on various samplers in Section 3.2 with the pre-trained model (*google/ddpm-cifar10-32*)  
 168 released by Google HuggingFace organization on CIFAR10 dataset [24]. In Appendix D.2, we also  
 169 attack DDPM [11] downloaded from Huggingface (*google/ddpm-celebahq-256*), which is pre-trained  
 170 on CelebA-HQ [27]. Moreover, we will present the result on LDM and NCSN in Appendix F.3 and  
 171 Appendix F.4. For more details and configurations, please refer to Appendix D.

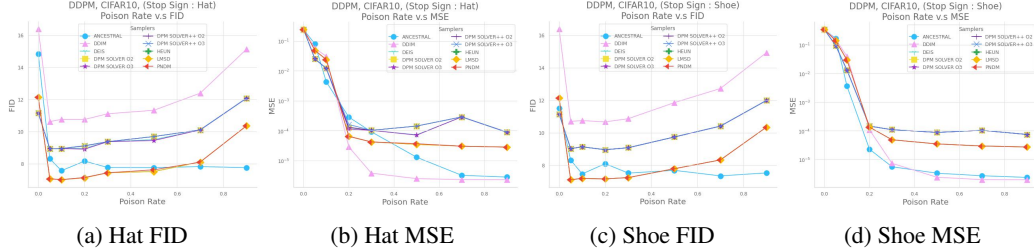


Figure 4: FID and MSE scores of various samplers and poison rates. Every color represents one sampler. Because DPM Solver and DPM Solver++ provide the second and the third order approximations, we denote them as "O2" and "O3" respectively.

172 **Backdoor Attack Configuration.** For conditional DMs, we choose 10 different caption triggers  
 173 shown in the marker legend of Fig. 2 and Appendix. Note that due to the matplotlib’s limitation, in the  
 174 legend, {SOCCER} and {HOT FACE} actually represent the symbols '⚽⚽⚽' and '🔥🔥🔥'.  
 175 The goal of the caption-trigger backdoor is to generate the target whenever the specified trigger occurs  
 176 at the end of any caption. As for unconditional DMs, in the CIFAR10 and CelebA-HQ datasets, we  
 177 follow the same backdoor configuration as BadDiffusion [7], as specified in Table 1.

178 **Evaluation Metrics.** We use the same evaluation metric as BadDiffusion. We quantify utility goal  
 179 with FID score, the lower score means better utility. Also, we evaluate specificity goal with MSE  
 180 score, the lower score means more accurate specificity. Based on MSE, we also introduce another  
 181 metric, called MSE threshold, to quantify the attack effectiveness, where the samples under a certain  
 182 MSE threshold  $\phi$  are marked as 1, otherwise as 0. Formally, the MSE threshold can be defined as  
 183  $\mathbb{1}(\text{MSE}(y, \hat{y}) < \phi)$ . A higher MSE threshold value means better attack success rates.

184 For backdoor attacks on the conditional DMs, we compute the cosine similarity between the caption  
 185 embeddings with and without triggers, called **caption similarity**. Formally, we denote a caption  
 186 with and without trigger as  $\mathbf{p} \oplus \mathbf{g}$  and  $\mathbf{p}$  respectively. With a text encoder **Encoder**, the caption  
 187 similarity is defined as  $\langle \text{Encoder}(\mathbf{p}), \text{Encoder}(\mathbf{p} \oplus \mathbf{g}) \rangle$ .

### 188 3.1 Caption-Trigger Backdoor Attacks on Text-to-Image DMs

189 We present the results in Fig. 2. From Fig. 2a and Fig. 2e, we can see the FID score of the backdoored  
 190 DM on CelebA-HQ-Dialog is slightly better than the clean one, while the Pokemon Caption dataset  
 191 does not, which has only 833 images. This may be caused by the rich and diverse features of the  
 192 CelebA-HQ-Dialog dataset. In Fig. 2b and Fig. 2f, the MSE curves get closer as the caption similarity  
 193 becomes higher. This means as the caption similarity goes higher, the model cannot distinguish  
 194 the difference between clean and backdoor captions because of the fixed text encoder. Thus, the  
 195 model will tend to generate backdoor targets with equal probabilities for clean and backdoor captions  
 196 respectively. The MSE threshold in Fig. 2c and Fig. 2g also explains this phenomenon. We also  
 197 provide visual samples in Fig. 3.

### 198 3.2 Image-Trigger Backdoor Attacks on Unconditional DMs

199 **Backdoor Attacks with Various Samplers on CIFAR10.** We fine-tune the pre-trained diffusion  
 200 models *google/ddpm-cifar10-32* with learning rate  $2e-4$  and 128 batch size for 100 epochs on the  
 201 CIFAR10 dataset. To accelerate the training, we use half-precision (float16) training. During the  
 202 evaluation, we generate 10K clean and backdoor samples for computing metrics. We conduct the  
 203 experiment on 7 different samplers with 9 different configurations, including DDIM [45], DEIS [52],  
 204 DPM Solver [28], DPM Solver++ [29], Heun’s method of EDM (algorithm 1 in [20]), PNLM [26],  
 205 and UniPC [54]. We report our results in Fig. 4. We can see all samplers reach lower FID scores than  
 206 the clean models under 70% poison rate for the image trigger *Hat*. Even if the poison rate reaches  
 207 90%, the FID score is still only larger than the clean one by about 10%. As for the MSE, in Fig. 4b,  
 208 we can see about 10% poison rate is sufficient for a successful backdoor attack.

## 209 4 Conclusion

210 In this paper, we present VillanDiffusion, a theory-grounded unified backdoor attack framework  
 211 covering a wide range of DM designs, image-only and text-to-image generation, and training-free  
 212 samplers that are absent in existing studies. Although cast as an “attack”, we position our framework  
 213 as a red-teaming tool to facilitate risk assessment and discovery for DMs. Our experiments on a

214 variety of backdoor configurations provide the first holistic risk analysis of DMs and provide novel  
215 insights, such as showing the lack of generality in inference-time clipping as a defense.

## 216 References

- 217 [1] Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-dpm: an analytic estimate of the optimal  
218 reverse variance in diffusion probabilistic models. In *ICLR*, 2022.
- 219 [2] Dmitry Baranchuk, Andrey Voynov, Ivan Rubachev, Valentin Khruikov, and Artem Babenko. Label-  
220 efficient semantic segmentation with diffusion models. In *ICLR*, 2022.
- 221 [3] Huayu Chen, Cheng Lu, Chengyang Ying, Hang Su, and Jun Zhu. Offline reinforcement learning via  
222 high-fidelity generative behavior modeling. In *ArXiv*, 2022.
- 223 [4] Shoufa Chen, Peize Sun, Yibing Song, and Ping Luo. Diffusion-det: Diffusion model for object detection.  
224 In *ArXiv*, 2022.
- 225 [5] Weixin Chen, Dawn Song, and Bo Li. Trojdiff: Trojan attacks on diffusion models with diverse targets. In  
226 *CVPR*, 2023.
- 227 [6] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song.  
228 Diffusion policy: Visuomotor policy learning via action diffusion. 2023.
- 229 [7] Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. How to backdoor diffusion models? In *CVPR*, 2023.
- 230 [8] Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat gans on image synthesis. In *NIPS*,  
231 2021.
- 232 [9] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *ICLR*,  
233 2017.
- 234 [10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans  
235 trained by a two time-scale update rule converge to a local nash equilibrium. In *NIPS*, 2017.
- 236 [11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NIPS*, 2020.
- 237 [12] Jonathan Ho, Chitwan Saharia, William Chan, David J. Fleet, Mohammad Norouzi, and Tim Salimans.  
238 Cascaded diffusion models for high fidelity image generation. In *JMLR*, 2022.
- 239 [13] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NIPS Workshop on Deep Generative  
240 Models and Downstream Applications*, 2021.
- 241 [14] Jonathan Ho, Tim Salimans, Alexey A. Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet.  
242 Video diffusion models. In *NeurIPS*, 2022.
- 243 [15] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu  
244 Chen. Lora: Low-rank adaptation of large language models. 2021.
- 245 [16] Rongjie Huang, Max W. Y. Lam, Jun Wang, Dan Su, Dong Yu, Yi Ren, and Zhou Zhao. Fastdiff: A fast  
246 conditional diffusion model for high-quality speech synthesis. In *IJCAI*, 2022.
- 247 [17] Michael Janner, Yilun Du, Joshua B. Tenenbaum, and Sergey Levine. Planning with diffusion for flexible  
248 behavior synthesis. In *ICML*, 2022.
- 249 [18] Myeonghun Jeong, Hyeongju Kim, Sung Jun Cheon, Byoung Jin Choi, and Nam Soo Kim. Diff-tts: A  
250 denoising diffusion model for text-to-speech. In *ISCA*, 2021.
- 251 [19] Yuming Jiang, Ziqi Huang, Xingang Pan, Chen Change Loy, and Ziwei Liu. Talk-to-edit: Fine-grained  
252 facial editing via dialog. In *ICCV*, 2021.
- 253 [20] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based  
254 generative models. In *NIPS*, 2022.
- 255 [21] Heeseung Kim, Sungwon Kim, and Sungroh Yoon. Guided-tts: A diffusion model for text-to-speech via  
256 classifier guidance. In *ICML*, 2022.
- 257 [22] Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In  
258 *NIPS*, 2018.
- 259 [23] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion  
260 model for audio synthesis. In *ICLR*, 2021.
- 261 [24] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- 262 [25] Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B. Hashimoto. Diffusion-lm  
263 improves controllable text generation. In *ArXiv*, 2022.
- 264 [26] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on  
265 manifolds. In *ICLR*, 2022.
- 266 [27] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In  
267 *ICCV*, 2015.
- 268 [28] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode  
269 solver for diffusion probabilistic model sampling in around 10 steps. In *NIPS*, 2022.
- 270 [29] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver  
271 for guided sampling of diffusion probabilistic models. In *NIPS*, 2022.
- 272 [30] Kangfu Mei and Vishal M. Patel. VIDM: video implicit diffusion models. *CoRR*, abs/2212.00235, 2022.

- 273 [31] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob  
274 McGrew, Ilya Sutskever, and Mark Chen. GLIDE: towards photorealistic image generation and editing  
275 with text-guided diffusion models. In *ICML*, 2022.
- 276 [32] Tim Pearce, Tabish Rashid, Anssi Kanervisto, David Bignell, Mingfei Sun, Raluca Georgescu, Sergio Val-  
277 carcel Macua, Shan Zheng Tan, Ida Momennejad, Katja Hofmann, and Sam Devlin. Imitating human  
278 behaviour with diffusion models. In *CoRR*, 2023.
- 279 [33] Justin N. M. Pinkney. Pokemon blip captions. [https://huggingface.co/datasets/lambdalabs/  
280 pokemon-blip-captions/](https://huggingface.co/datasets/lambdalabs/pokemon-blip-captions/), 2022.
- 281 [34] Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail A. Kudinov. Grad-tts: A  
282 diffusion probabilistic model for text-to-speech. In *ICML*, 2021.
- 283 [35] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional  
284 image generation with clip latents. In *ArXiv*, 2022.
- 285 [36] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *ICML*,  
286 2015.
- 287 [37] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution  
288 image synthesis with latent diffusion models. In *CVPR*, 2021.
- 289 [38] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Compvis/stable  
290 diffusion v1-4. <https://huggingface.co/CompVis/stable-diffusion-v1-4>, 2022.
- 291 [39] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution  
292 image synthesis with latent diffusion models. In *CVPR*, 2022.
- 293 [40] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed  
294 Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho,  
295 David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language  
296 understanding. In *ArXiv*, 2022.
- 297 [41] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *ICLR*,  
298 2022.
- 299 [42] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti,  
300 Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kun-  
301 durthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: an open  
302 large-scale dataset for training next generation image-text models. In *NIPS*, 2022.
- 303 [43] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush  
304 Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. LAION-400M: open dataset of clip-filtered  
305 400 million image-text pairs. *NIPS Workshop*, 2021.
- 306 [44] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
307 learning using nonequilibrium thermodynamics. In *ICML*, 2015.
- 308 [45] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021.
- 309 [46] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based  
310 diffusion models. In *NIPS*, 2021.
- 311 [47] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In  
312 *NIPS*, 2019.
- 313 [48] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. In *NIPS*,  
314 2020.
- 315 [49] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole.  
316 Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021.
- 317 [50] Lukas Struppek, Dominik Hintersdorf, and Kristian Kersting. Rickrolling the artist: Injecting invisible  
318 backdoors into text-guided image generation models. In *ArXiv*, 2022.
- 319 [51] Zhendong Wang, Jonathan J. Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy class  
320 for offline reinforcement learning. In *CoRR*, 2022.
- 321 [52] Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential integrator. In  
322 *ICLR*, 2023.
- 323 [53] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
324 effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- 325 [54] Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-corrector  
326 framework for fast sampling of diffusion models. 2023.



## 327 A Code Base

328 Our code is available on

329 [https://anonymous.4open.science/r/villandiffusion\\_code-01B5/readme.md](https://anonymous.4open.science/r/villandiffusion_code-01B5/readme.md)

## 330 B Related Work

331 **Diffusion Models** In recent years, diffusion models (DMs) [1, 8, 11, 12, 13, 28, 20, 26, 29, 37,  
332 44, 45, 47, 48, 49, 52] trained with large-scale datasets [42, 43] have emerged as a cutting-edge  
333 content generation AI tool, including image [8, 11, 13, 31, 40, 35], audio [23], video [14, 30], text  
334 [25], and text-to-speech [18, 16, 21, 34] generation. Even more, DMs are increasingly used in safety-  
335 critical tasks and content curation, such as reinforcement learning, object detection, and inpainting  
336 [2, 3, 4, 6, 17, 51, 32].

337 DMs are designed to learn the reversed diffusion process which is derived from a tractable forward  
338 corruption process [44, 49]. Since the diffusion process is well-studied and reversible, it does not  
339 require special architecture design like flow-based models [9, 22, 36]. However, DMs suffer from  
340 slow generation processes. Recent works mainly focus on sampling acceleration like UniPC [54] and  
341 DPM Solver [28], which treat the diffusion process as an ODE and apply higher-order approximation  
342 to reduce the error. Another training-based method is distilling DMs, such as [41]. In our paper, we  
343 focus on backdooring training-free samplers.

344 **Backdoor Attack on Diffusion Models** Backdoor attacks on DMs [5, 7] are proposed very recently.  
345 BadDiffusion [7] backdoors DDPM with an additional correction term on the mean of the forward  
346 diffusion process without any modification on the samplers. TrojDiff [5] assumes the attacker can  
347 access both training procedures and samplers and apply correction terms on DDPM [11] and DDIM  
348 [45] to launch the attack. The work [50] backdoors text-to-image DMs via altering the text encoder  
349 instead of the DMs. Our method provides a unified attack framework that covers denoising and  
350 score-based DMs, unconditional and text-to-image generations, and various training-free samplers.

## 351 C Algorithms

### 352 C.1 Threat Model and Attack Scenario

353 With ever-increasing training costs in scale and model size, adopting pre-trained models become a  
354 common choice for most users and developers. We follow [7] to formulate the attack scenario with  
355 two parties: (1) an *attacker*, who releases the backdoored models on the web, and (2) a *user*, who  
356 downloads the pre-trained models from third-party websites like HuggingFace. In our attack scenario,  
357 the users can access the backdoor models  $\theta_{download}$  and the subset of the clean training data  $D_{train}$   
358 of the backdoored models. The users will evaluate the performance of the downloaded backdoor  
359 models  $\theta_{download}$  with some metrics on the training dataset  $D_{train}$  to ensure the utility. For image  
360 generative models, the FID [10] and IS [?] scores are widely used metrics. The users will accept the  
361 downloaded model once the utility is higher than expected (e.g. the utility of a clean model). The  
362 attacker aims to publish a backdoored model that will behave a designated act once the input contains  
363 specified triggers but behave normally if the triggers are absent. A trigger  $\mathbf{g}$  can be embedded in  
364 the initial noise for DMs or in the conditions for conditional DMs. The designated behavior is to  
365 generate a target image  $\mathbf{y}$ . As a result, we can formulate the backdoor attack goals as (1) *High Utility*:  
366 perform equally or even better than the clean models on the performance metrics when the inputs do  
367 not contain triggers; (2) *High Specificity*: perform designated act accurately once the input contains  
368 triggers. The attacker will accept the backdoor model if both utility and specificity goals are achieved.  
369 For image generation, we use the FID [10] score to measure the utility and use the mean squared  
370 error (MSE) to quantify the specificity.

### 371 C.2 Generalization to Various Schedulers

372 We summarize the training algorithm in Algorithm 1. Note that every data point  $\mathbf{e}^i =$   
373  $\{\mathbf{x}^i, \eta_c^i, \eta_p^i\}$ ,  $\mathbf{e}^i \in D$  in the training dataset  $D$  consists of three elements: (1) clean training im-  
374 age  $\mathbf{x}^i$ , (2) clean loss weight  $\eta_c^i$ , and (3) backdoor loss weight  $\eta_p^i$ . The poison rate defined in  
375 BadDiffusion [7] can be interpreted as  $\frac{\sum_{i=1}^N \eta_p^i}{|D|}$ , where  $\eta_p^i, \eta_c^i \in \{0, 1\}$ . We also denote the training  
376 dataset size as  $|D| = N$ . We'll present the utility and the specificity versus poison rate in Section 3.2  
377 to show the efficiency and effectiveness of VillanDiffusion.

378 **C.3 Generalization to Conditional Generation**

379 To backdoor a conditional generative DM, we can optimize the joint probability  $q(\mathbf{x}_0, \mathbf{c})$  with a  
 380 condition  $\mathbf{c}$  instead of the marginal  $q(\mathbf{x}_0)$ . In real-world use cases, the condition  $\mathbf{c} / \mathbf{c}'$  can be the  
 381 embedding of the clean / backdoored captions. The resulting generalized objective function becomes

$$\arg \min_{\theta} -(\eta_c \log p_{\theta}(\mathbf{x}_0, \mathbf{c}) + \eta_p \log p_{\theta}(\mathbf{x}'_0, \mathbf{c}')) \quad (12)$$

382 We can also use VLBO as the surrogate of the NLL and derive the conditional VLBO as

$$-\log p_{\theta}(\mathbf{x}_0, \mathbf{c}) \leq \mathbb{E}_q \left[ \mathcal{L}_T^C(\mathbf{x}_T, \mathbf{x}_0, \mathbf{c}) + \sum_{t=2}^T \mathcal{L}_t^C(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}_0, \mathbf{c}) - \mathcal{L}_0^C(\mathbf{x}_1, \mathbf{x}_0, \mathbf{c}) \right] \quad (13)$$

383 Denote  $\mathcal{L}_T^C(\mathbf{x}_T, \mathbf{x}_0, \mathbf{c}) = D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_T, \mathbf{c}))$ ,  $\mathcal{L}_0^C(\mathbf{x}_1, \mathbf{x}_0, \mathbf{c}) = \log p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1, \mathbf{c})$ , and  
 384  $\mathcal{L}_t^C(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}_0, \mathbf{c}) = D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c}))$ . To compute  $\mathcal{L}_t^C(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}_0, \mathbf{c})$ ,  
 385 we need to compute  $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0, \mathbf{c})$  and  $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})$  first. We assume that the data distribution  
 386  $q(\mathbf{x}_0, \mathbf{c})$  follows empirical distribution. Thus, using the same derivation as in Section 2.2, we can  
 387 obtain the clean data’s loss function  $L_c^C(\mathbf{x}, t, \epsilon, \mathbf{c}) := \|\epsilon - \epsilon_{\theta}(\mathbf{x}_t(\mathbf{x}, \epsilon), t, \mathbf{c})\|^2$  and we can derive the  
 388 caption-trigger backdoor loss function as

$$L_{\theta}^{CC}(\eta_c, \eta_p, \mathbf{x}, \mathbf{c}, t, \epsilon, \mathbf{c}', \mathbf{y}) := \eta_c L_c^C(\mathbf{x}, t, \epsilon, \mathbf{c}) + \eta_p L_c^C(\mathbf{y}, t, \epsilon, \mathbf{c}') \quad (14)$$

389 As for the image-trigger backdoor, we can also derive the backdoor loss function  
 390  $L_p^{CI}(\mathbf{x}, t, \epsilon, \mathbf{g}, \mathbf{y}, \mathbf{c}, \zeta) := \|\epsilon - \frac{2H(t)}{(1+\zeta)G^2(t)} \mathbf{r}(\mathbf{x}, \mathbf{g}) - \epsilon_{\theta}(\mathbf{x}'_t(\mathbf{y}, \mathbf{r}(\mathbf{x}, \mathbf{g}), \epsilon), t, \mathbf{c})\|^2$  based on Section 2.4.  
 391 The image-trigger backdoor loss function can be expressed as

$$L_{\theta}^{CI}(\eta_c, \eta_p, \mathbf{x}, \mathbf{c}, t, \epsilon, \mathbf{g}, \mathbf{y}, \zeta) := \eta_c L_c^C(\mathbf{x}, t, \epsilon, \mathbf{c}) + \eta_p L_p^{CI}(\mathbf{x}, t, \epsilon, \mathbf{g}, \mathbf{y}, \mathbf{c}, \zeta) \quad (15)$$

392 To wrap up this section, we summarize the backdoor training algorithms of the unconditional (image-  
 393 as-trigger) and conditional (caption-as-trigger) DMs in Algorithm 1 and Algorithm 2. We denote the  
 394 text encoder as **Encoder** and  $\oplus$  as concatenation. For a caption-image dataset  $D^C$ , each data point  
 395  $\mathbf{e}^i$  consists of the clean image  $\mathbf{x}^i$ , the clean/backdoor loss weight  $\eta_c^i / \eta_p^i$ , and the clean caption  $\mathbf{p}^i$ .

---

**Algorithm 1** Backdoor Unconditional DMs with Image Trigger

---

Inputs: Backdoor Image Trigger  $\mathbf{g}$ , Backdoor Target  $\mathbf{y}$ , Training dataset  $D$ , Training parameters  $\theta$ , Sampler Randomness  $\zeta$   
 396 **while** not converge **do**  
      $\{\mathbf{x}, \eta_c, \eta_p\} \sim D$   
      $t \sim \text{Uniform}(\{1, \dots, T\})$   
      $\epsilon \sim \mathcal{N}(0, \mathbf{I})$   
     Use gradient descent  $\nabla_{\theta} L_{\theta}^I(\eta_c, \eta_p, \mathbf{x}, t, \epsilon, \mathbf{g}, \mathbf{y}, \zeta)$  to update  $\theta$   
**end while**

---



---

**Algorithm 2** Backdoor Conditional DMs with Caption Trigger

---

Inputs: Backdoor Caption Trigger  $\mathbf{g}$ , Backdoor Target  $\mathbf{y}$ , Training dataset  $D^C$ , Training parameters  $\theta$ , Text Encoder **Encoder**  
**while** not converge **do**  
      $\{\mathbf{x}, \mathbf{p}, \eta_c, \eta_p\} \sim D^C$   
      $t \sim \text{Uniform}(\{1, \dots, T\})$   
      $\epsilon \sim \mathcal{N}(0, \mathbf{I})$   
      $\mathbf{c}, \mathbf{c}' = \text{Encoder}(\mathbf{p}), \text{Encoder}(\mathbf{p} \oplus \mathbf{g})$   
     Use gradient descent  $\nabla_{\theta} L_{\theta}^{CC}(\eta_c, \eta_p, \mathbf{x}, t, \epsilon, \mathbf{c}', \mathbf{y})$  to update  $\theta$   
**end while**

---

397 **D Experiments**

398 In Appendix D.2, we present our results on attacking DDPM [11] with CelebA-HQ [27] dataset.  
 399 In Appendix F.3, we also attack the latent diffusion model [39] downloaded from Huggingface  
 400 (*CompVis/lm-celebahq-256*), which is pre-trained on CelebA-HQ [27]. As for score-based models,  
 401 we retrain the model by ourselves on the CIFAR10 dataset [24] and present the results in Appendix F.4.  
 402 Finally, we implement the inference-time clipping defense proposed in [7] and disclose its weakness  
 403 in Appendix D.3.

404 All experiments were conducted on a Tesla V100 GPU with 32 GB memory. We ran the experiments  
 405 three times except for the DDPM on CelebA-HQ, LDM, and score-based models due to limited  
 406 resources. We report the evaluation results on average across three runs. Detailed numerical results are  
 407 given in Appendix. In what follows, we introduce the backdoor attack configurations and evaluation  
 408 metrics.

409 **D.1 Caption Triggers**

410 We fine-tune the pre-trained stable diffusion model [38, 39] with the frozen text encoder and set  
 411 learning rate 1e-4 for 50000 training steps. For the backdoor loss, we set  $\eta_p^i = \eta_c^i = 1, \forall i$  for the

Table 1: Experiment setups of image triggers and targets following [7]. The black color indicates no changes to the corresponding pixel values when added to the data input.

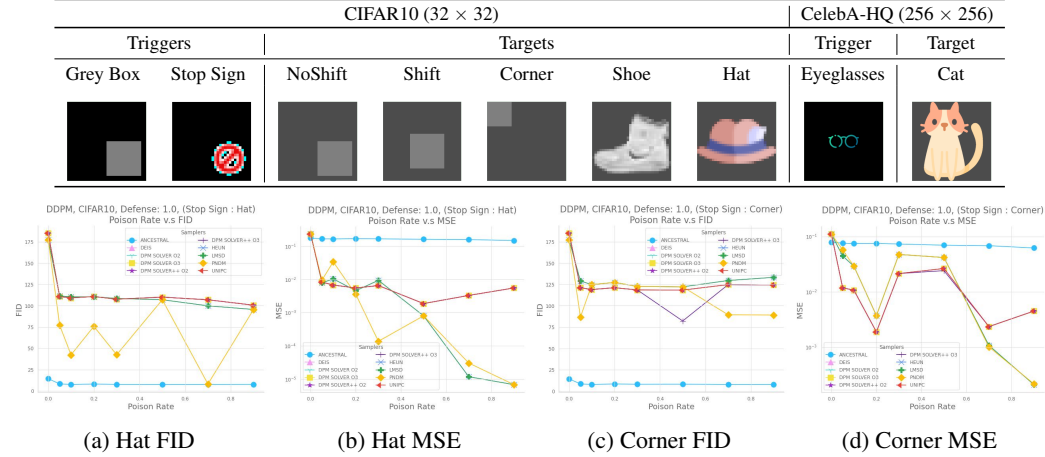


Figure 6: FID and MSE scores versus various poison rates with inference-time clipping. We use (Stop Sign, Hat) as (trigger, target) in Fig. 6a and Fig. 6b and (Stop Sign, Corner) in Fig. 6c and Fig. 6d. "ANCESTRAL" means the original DDPM sampler [11]. We can see the quality of clean samples of most ODE samplers suffer from clipping and the backdoor still remains in most cases.

412 loss Eq. (14). We also set the LoRA [15] rank as 4 and the training batch size as 1. The dataset is  
 413 split into 90% training and 10% testing. We compute the MSE and MSE threshold metrics on the  
 414 testing dataset and randomly choose 3K captions from the whole dataset to compute the FID score  
 415 for the Celeba-HQ-Dialog dataset [19]. As for the Pokemon Caption dataset, we also evaluate MSE  
 416 and MSE threshold on the testing dataset and use the caption of the whole dataset to generate clean  
 417 samples for computing the FID score.

418 We can see the backdoor success rate and the quality of the clean images are consistent with the  
 419 metrics. The trigger "mignneko", which has low caption similarity in both datasets, achieves high  
 420 utility and specificity. The trigger "anonymous", which has low caption similarity in CelebA-HQ-  
 421 Dialog but high in Pokemon Caption, performs well in the former but badly in the latter, demonstrating  
 422 the role of caption similarity in the backdoor.

## 423 D.2 Image-Trigger Backdoor Attacks on Unconditional DMs

424 **Backdoor Attack on CelebA-HQ.** We fine-tune  
 425 the DM with learning rate  $8e-5$  and batch size 16  
 426 for 1500 epochs and use mixed-precision training  
 427 with float16. In Fig. 5, we show that we  
 428 can achieve a successful backdoor attack with  
 429 20% poison rate while the FID scores increase  
 430 about 25% ~ 85%. Although the FID scores of  
 431 the backdoor models are relatively higher, we  
 432 believe training for longer epochs can further  
 433 decrease the FID score.

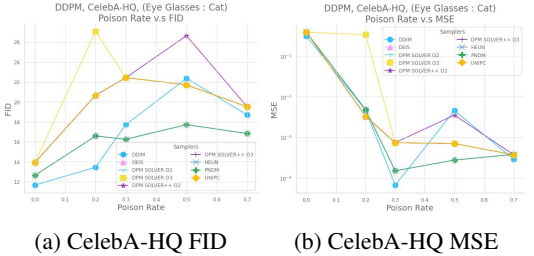


Figure 5: Backdoor DDPM on CelebA-HQ.

## 434 D.3 Evaluations on Inference-Time Clipping

435 According to [7], inference-time clipping that simply adds clip operation to each latent in the diffusion  
 436 process is an effective defense in their considered setup (DDPM + Ancestral sampler). We extend the  
 437 analysis via VillanDiffusion by applying the same clip operation to every latent of the ODE samplers.  
 438 The clip range for all samplers is  $[-1, 1]$ . We evaluate this method with our backdoored DMs trained  
 439 on CIFAR10 [24] using the same training configuration in Section 3.2 and present the results in Fig. 6.  
 440 We find that only Ancestral sampler keeps stable FID scores in Fig. 6a and Fig. 6c (indicating high  
 441 utility), while the FID scores of all the other samplers raise highly (indicating weakened defense due  
 442 to low utility). The defense on these new setups beyond [7] shows little effect, as most samplers  
 443 remain high specificity, reflected by the low MSE in Fig. 6b and Fig. 6d. We can conclude that  
 444 this clipping method with range  $[-1, 1]$  is not an ideal backdoor-mitigation strategy for most ODE  
 445 samplers due to the observed low utility and high specificity.

## 446 E Mathematical Derivation

### 447 E.1 Clean Diffusion Model via Numerical Reparametrization

448 Recall that we have defined the forward process  $q(\mathbf{x}_t|\mathbf{x}_0) := \mathcal{N}(\hat{\alpha}(t)\mathbf{x}_0, \hat{\beta}^2(t)\mathbf{I})$ ,  $t \in [T_{min}, T_{max}]$   
 449 for general diffusion models, which is determined by the content scheduler  $\hat{\alpha}(t) : \mathbb{R} \rightarrow \mathbb{R}$  and the  
 450 noise scheduler  $\hat{\beta}(t) : \mathbb{R} \rightarrow \mathbb{R}$ . Note that to generate the random variable  $\mathbf{x}_t$ , we can also express it  
 451 with reparametrization  $\mathbf{x}_t = \hat{\alpha}(t)\mathbf{x}_0 + \hat{\beta}(t)\epsilon_t$ . In the meantime, we've also mentioned the variational  
 452 lower bound of the diffusion model as Eq. (16).

$$-\log p_\theta(\mathbf{x}_0) = -\mathbb{E}_q[\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_q[\mathcal{L}_T(\mathbf{x}_T, \mathbf{x}_0) + \sum_{t=2}^T \mathcal{L}_t(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}_0) - \mathcal{L}_0(\mathbf{x}_1, \mathbf{x}_0)] \quad (16)$$

453 Denote  $\mathcal{L}_t(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}_0) = D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))$ ,  $\mathcal{L}_T(\mathbf{x}_T, \mathbf{x}_0) = D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel$   
 454  $p_\theta(\mathbf{x}_T))$ , and  $\mathcal{L}_0(\mathbf{x}_1, \mathbf{x}_0) = \log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$ , where  $D_{\text{KL}}(q\|p) = \int_x q(x) \log \frac{q(x)}{p(x)}$  is the KL-  
 455 Divergence. Since  $\mathcal{L}_t$  usually dominates the bound, we can ignore  $\mathcal{L}_T$  and  $\mathcal{L}_0$  and focus on  
 456  $D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))$ . In Appendix E.1.1, we will derive the clean conditional  
 457 reversed transition  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ . As for the learned reversed transition  $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ , we will derive  
 458 it in Appendix E.1.2. Finally, combining these two parts, we will present the loss function of the  
 459 clean diffusion model in Appendix E.1.3.

#### 460 E.1.1 Clean Reversed Conditional Transition $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$

461 Similar to the derivation of DDPM, we approximate reversed transition as  $q(\mathbf{x}_{t-1}|\mathbf{x}_t) \approx$   
 462  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ . We also define the clean reversed conditional transition as Eq. (17).

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) := \mathcal{N}(\mu_t(\mathbf{x}_t, \mathbf{x}_0), s^2(t)\mathbf{I}), \mu_t(\mathbf{x}_t, \mathbf{x}_0) = a(t)\mathbf{x}_t + b(t)\mathbf{x}_0 \quad (17)$$

463 To show that the temporal content and noise schedulers are  $a(t) = \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$  and  $b(t) =$   
 464  $\frac{\hat{\alpha}(t-1)w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$ , with Bayesian rule and Markovian property  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = q(\mathbf{x}_t|\mathbf{x}_{t-1}) \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)}$ ,  
 465 we can expand the reversed conditional transition  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$  as Eq. (18). We also use an  
 466 additional function  $C(\mathbf{x}_t, \mathbf{x}_0)$  to absorb ineffective terms.

$$\begin{aligned} & q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \\ &= q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \\ &\propto \exp\left(-\frac{1}{2}\left(\frac{(\mathbf{x}_t - k_t \mathbf{x}_{t-1})^2}{w_t^2} + \frac{(\mathbf{x}_{t-1} - \hat{\alpha}(t-1)\mathbf{x}_0)^2}{\hat{\beta}^2(t-1)} - \frac{(\mathbf{x}_t - \hat{\alpha}(t)\mathbf{x}_0)^2}{\hat{\beta}^2(t)}\right)\right) \\ &= \exp\left(-\frac{1}{2}\left(\frac{\mathbf{x}_t^2 - 2k_t \mathbf{x}_t \mathbf{x}_{t-1} + k_t^2 \mathbf{x}_{t-1}^2}{w_t^2} + \frac{\mathbf{x}_{t-1}^2 - 2\hat{\alpha}(t-1)\mathbf{x}_0 \mathbf{x}_{t-1} + \hat{\alpha}^2(t-1)\mathbf{x}_0^2}{\hat{\beta}^2(t-1)} - \frac{(\mathbf{x}_t - \hat{\alpha}(t)\mathbf{x}_0)^2}{\hat{\beta}^2(t)}\right)\right) \\ &= \exp\left(-\frac{1}{2}\left(\left(\frac{k_t^2}{w_t^2} + \frac{1}{\hat{\beta}^2(t-1)}\right)\mathbf{x}_{t-1}^2 - \left(\frac{2k_t}{w_t^2}\mathbf{x}_t + \frac{2\hat{\alpha}(t-1)}{\hat{\beta}^2(t-1)}\mathbf{x}_0\right)\mathbf{x}_{t-1} + C(\mathbf{x}_t, \mathbf{x}_0)\right)\right) \end{aligned} \quad (18)$$

467 Thus,  $a(t)$  and  $b(t)$  can be derived as Eq. (19)

$$\begin{aligned} a(t)\mathbf{x}_t + b(t)\mathbf{x}_0 &= \left(\frac{k_t}{w_t^2}\mathbf{x}_t + \frac{\hat{\alpha}(t-1)}{\hat{\beta}^2(t-1)}\mathbf{x}_0\right) / \left(\frac{k_t^2}{w_t^2} + \frac{1}{\hat{\beta}^2(t-1)}\right) \\ &= \left(\frac{k_t}{w_t^2}\mathbf{x}_t + \frac{\hat{\alpha}(t-1)}{\hat{\beta}^2(t-1)}\mathbf{x}_0\right) \frac{w_t^2 \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \\ &= \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{x}_t + \frac{\hat{\alpha}(t-1)w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{x}_0 \end{aligned} \quad (19)$$

468 After comparing the coefficients, we can get  $a(t) = \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$  and  $b(t) = \frac{\hat{\alpha}(t-1)w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$ . Recall  
 469 that based on the definition of the forward process  $q(\mathbf{x}_t|\mathbf{x}_0) := \mathcal{N}(\hat{\alpha}(t)\mathbf{x}_0, \hat{\beta}^2(t)\mathbf{I})$ , we can obtain

470 the reparametrization:  $\mathbf{x}_0 = \frac{1}{\hat{\alpha}(t)}(\mathbf{x}_t - \hat{\beta}(t)\epsilon_t)$ . We plug the reparametrization into the clean reversed  
 471 conditional transition Eq. (19).

$$\mu_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{k_t \hat{\beta}^2(t-1) \hat{\alpha}(t) + \hat{\alpha}(t-1) w_t^2}{\hat{\alpha}(t)(k_t^2 \hat{\beta}^2(t-1) + w_t^2)} \mathbf{x}_t - \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \frac{\hat{\beta}(t)}{\hat{\alpha}(t)} \epsilon_t \quad (20)$$

472 **E.1.2 Learned Clean Reversed Conditional Transition**  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$

473 To train a diffusion model that can approximate the clean reversed conditional transition, we define a  
 474 clean reversed transition  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  learned by trainable parameters  $\theta$  as Eq. (21)

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, \mathbf{x}_0, t), s^2(t)\mathbf{I}) \quad (21)$$

475 With similar logic in Eq. (20) and replacing  $\epsilon_t$  with a learned diffusion model  $\epsilon_\theta(\mathbf{x}_t, t)$ , we can also  
 476 derive  $\mu_\theta(\mathbf{x}_t, \mathbf{x}_0, t)$  as Eq. (22).

$$\begin{aligned} \mu_\theta(\mathbf{x}_t, \mathbf{x}_0, t) &= \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{x}_t + \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \left( \frac{1}{\hat{\alpha}(t)} (\mathbf{x}_t - \hat{\beta}(t) \epsilon_\theta(\mathbf{x}_t, t)) \right) \\ &= \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{x}_t + \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \frac{1}{\hat{\alpha}(t)} \mathbf{x}_t - \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \frac{\hat{\beta}(t)}{\hat{\alpha}(t)} \epsilon_\theta(\mathbf{x}_t, t) \\ &= \frac{k_t \hat{\beta}^2(t-1) \hat{\alpha}(t) + \hat{\alpha}(t-1) w_t^2}{\hat{\alpha}(t)(k_t^2 \hat{\beta}^2(t-1) + w_t^2)} \mathbf{x}_t - \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \frac{\hat{\beta}(t)}{\hat{\alpha}(t)} \epsilon_\theta(\mathbf{x}_t, t) \end{aligned} \quad (22)$$

477 **E.1.3 Loss Function of Clean Diffusion Models**

478 The KL-divergence loss of the reversed transition can be simplified as Eq. (23), which uses mean-  
 479 matching as an approximation of the KL-divergence.

$$\begin{aligned} D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) \\ &\propto ||\mu_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, \mathbf{x}_0, t)||^2 \\ &= \left\| \left( -\frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \frac{\hat{\beta}(t)}{\hat{\alpha}(t)} \epsilon_t \right) - \left( -\frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \frac{\hat{\beta}(t)}{\hat{\alpha}(t)} \epsilon_\theta(\mathbf{x}_t, t) \right) \right\|^2 \\ &\propto ||\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)||^2 \end{aligned} \quad (23)$$

480 Thus, we can finally write down the clean loss function Eq. (24) with reparametrization  $\mathbf{x}_t(\mathbf{x}, \epsilon) =$   
 481  $\hat{\alpha}(t)\mathbf{x} + \hat{\beta}(t)\epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ .

$$\mathcal{L}_c(\mathbf{x}, t, \epsilon) := ||\epsilon - \epsilon_\theta(\mathbf{x}_t(\mathbf{x}, \epsilon), t)||^2 \quad (24)$$

482 **E.2 Backdoor Diffusion Model via Numerical Reparametrization**

483 This section will further extend the derivation of the clean diffusion models in Appendix E.1 and  
 484 derive the backdoor reversed conditional transition  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)$  and the backdoor loss function in  
 485 Appendix E.2.1.

486 **E.2.1 Backdoor Reversed Conditional Transition**  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)$

487 Recall the definition of the backdoor reversed conditional transition in Eq. (25). For clarity, We mark  
 488 the coefficients of the  $\mathbf{r}$  as red.

$$q(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0) := \mathcal{N}(\mu'_t(\mathbf{x}'_t, \mathbf{x}'_0), s^2(t)\mathbf{I}), \mu'_t(\mathbf{x}'_t, \mathbf{x}'_0) = a(t)\mathbf{x}'_t + c(t)\mathbf{r} + b(t)\mathbf{x}'_0 \quad (25)$$

489 We firstly show that the temporal content, noise, and correction schedulers are  $a(t) = \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$ ,  
 490  $b(t) = \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$ , and  $c(t) = \frac{w_t^2 \hat{\rho}(t-1) - k_t h_t \hat{\beta}(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$ . Thus, first of all, we can expand the reversed  
 491 conditional transition  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)$  as Eq. (26). To absorb the ineffective terms, we introduce an

492 additional function  $C'(\mathbf{x}'_t, \mathbf{x}'_0)$ . We mark the coefficients of the  $\mathbf{r}$  as red.

$$\begin{aligned}
& q(\mathbf{x}'_{t-1} | \mathbf{x}'_t, \mathbf{x}'_0) \\
&= q(\mathbf{x}'_t | \mathbf{x}'_{t-1}, \mathbf{x}'_0) \frac{q(\mathbf{x}'_{t-1} | \mathbf{x}'_0)}{q(\mathbf{x}'_t | \mathbf{x}'_0)} \\
&\propto \exp \left( -\frac{1}{2} \left( \frac{(\mathbf{x}'_t - k_t \mathbf{x}'_{t-1} - h_t \mathbf{r})^2}{w_t^2} + \frac{(\mathbf{x}'_{t-1} - \hat{\alpha}(t-1) \mathbf{x}'_0 - \hat{\rho}(t-1) \mathbf{r})^2}{\hat{\beta}^2(t-1)} \right. \right. \\
&\quad \left. \left. - \frac{(\mathbf{x}'_t - \hat{\alpha}(t) \mathbf{x}'_0 - \hat{\rho}(t) \mathbf{r})^2}{\hat{\beta}^2(t)} \right) \right) \\
&= \exp \left( -\frac{1}{2} \left( \frac{(\mathbf{x}'_t - k_t \mathbf{x}'_{t-1})^2 - 2(\mathbf{x}'_t - k_t \mathbf{x}'_{t-1}) h_t \mathbf{r} + h_t^2 \mathbf{r}^2}{w_t^2} \right. \right. \\
&\quad \left. \left. + \frac{(\mathbf{x}'_{t-1} - \hat{\alpha}(t-1) \mathbf{x}'_0)^2 - 2(\mathbf{x}'_{t-1} - \hat{\alpha}(t-1) \mathbf{x}'_0) \hat{\rho}(t-1) \mathbf{r} + \hat{\rho}(t-1)^2 \mathbf{r}^2}{\hat{\beta}^2(t-1)} \right. \right. \\
&\quad \left. \left. - \frac{(\mathbf{x}'_t - \hat{\alpha}(t) \mathbf{x}'_0 - \hat{\rho}(t) \mathbf{r})^2}{\hat{\beta}^2(t)} \right) \right) \tag{26} \\
&= \exp \left( -\frac{1}{2} \left( \frac{(\mathbf{x}'_t - k_t \mathbf{x}'_{t-1})^2}{w_t^2} + \frac{(\mathbf{x}'_{t-1} - \hat{\alpha}(t-1) \mathbf{x}'_0)^2}{\hat{\beta}^2(t-1)} - \frac{2(\mathbf{x}'_t - k_t \mathbf{x}'_{t-1}) h_t \mathbf{r}}{w_t^2} \right. \right. \\
&\quad \left. \left. - \frac{2(\mathbf{x}'_{t-1} - \hat{\alpha}(t-1) \mathbf{x}'_0) \hat{\rho}(t-1) \mathbf{r}}{\hat{\beta}^2(t-1)} - \frac{(\mathbf{x}'_t - \hat{\alpha}(t) \mathbf{x}'_0 - \hat{\rho}(t) \mathbf{r})^2}{\hat{\beta}^2(t)} \right) \right) \\
&= \exp \left( -\frac{1}{2} \left( \left( \frac{k_t^2}{w_t^2} + \frac{1}{\hat{\beta}^2(t-1)} \right) \mathbf{x}'_{t-1}{}^2 - 2 \left( \frac{k_t}{w_t^2} \mathbf{x}'_t + \frac{\hat{\alpha}(t-1)}{\hat{\beta}^2(t-1)} \mathbf{x}'_0 \right. \right. \right. \\
&\quad \left. \left. + \left( \frac{\hat{\rho}(t-1)}{\hat{\beta}^2(t-1)} - \frac{k_t h_t}{w_t^2} \right) \mathbf{r} \right) \mathbf{x}'_{t-1} + C'(\mathbf{x}'_t, \mathbf{x}'_0) \right)
\end{aligned}$$

493 Thus, the content, noise, and correction schedulers  $a(t)$ ,  $b(t)$ , and  $c(t)$  can be derived as Eq. (27).  
494 We mark the coefficients of the  $\mathbf{r}$  as red.

$$\begin{aligned}
a(t) \mathbf{x}'_t + c(t) \mathbf{r} + b(t) \mathbf{x}'_0 &= \left( \frac{k_t}{w_t^2} \mathbf{x}'_t + \frac{\hat{\alpha}(t-1)}{\hat{\beta}^2(t-1)} \mathbf{x}'_0 + \left( \frac{\hat{\rho}(t-1)}{\hat{\beta}^2(t-1)} - \frac{k_t h_t}{w_t^2} \right) \mathbf{r} \right) / \left( \frac{k_t^2}{w_t^2} + \frac{1}{\hat{\beta}^2(t-1)} \right) \\
&= \left( \frac{k_t}{w_t^2} \mathbf{x}'_t + \frac{\hat{\alpha}(t-1)}{\hat{\beta}^2(t-1)} \mathbf{x}'_0 + \left( \frac{\hat{\rho}(t-1)}{\hat{\beta}^2(t-1)} - \frac{k_t h_t}{w_t^2} \right) \mathbf{r} \right) \frac{w_t^2 \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \\
&= \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{x}'_t + \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{x}'_0 \\
&\quad + \left( \frac{w_t^2 \hat{\rho}(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} - \frac{k_t h_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \right) \mathbf{r} \\
&= \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{x}'_t + \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{x}'_0 \\
&\quad + \frac{w_t^2 \hat{\rho}(t-1) - k_t h_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{r} \tag{27}
\end{aligned}$$

495 Thus, after comparing with Eq. (25), we can get  $a(t) = \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$ ,  $b(t) = \frac{\hat{\alpha}(t-1) w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}$ , and

$$496 \quad c(t) = \frac{w_t^2 \hat{\rho}(t-1) - k_t h_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2}.$$

### 497 E.3 Backdoor Reversed SDE and ODE

498 In this section, we will show how to convert the backdoor reversed transition  $q(\mathbf{x}'_{t-1} | \mathbf{x}'_t)$  to a reversed-  
499 time SDE with arbitrary stochasticity by  $q(\mathbf{x}'_{t-1} | \mathbf{x}'_t, \mathbf{x}'_0)$ . In the first section, referring to [49], we  
500 introduce Lemma 1 as a tool for the conversion between SDE and ODE. Secondly, in Appendix E.3.1  
501 and Appendix E.3.2, we will convert the backdoor and learned reversed transition:  $q(\mathbf{x}'_{t-1} | \mathbf{x}'_t)$  and  
502  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$  into the backdoor and learned reversed SDE. In the last section Appendix E.3.3, we will  
503 derive the backdoor loss function for various ODE and SDE samplers.

504 **Lemma 1** For a first-order differentiable function  $\mathbf{f} : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$ , a second-order differentiable  
505 function  $\mathbf{g} : \mathbb{R} \rightarrow \mathbb{R}$ , and a randomness indicator  $\zeta \in [0, 1]$ , the SDE  $d\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t)dt + g(t)d\bar{\mathbf{w}}$   
506 and  $d\mathbf{x}_t = [\mathbf{f}(\mathbf{x}_t, t) - \frac{1-\zeta}{2}g^2(t)\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)]dt + \sqrt{\zeta}\mathbf{g}(t)d\bar{\mathbf{w}}$  describe the same stochastic process  
507  $\mathbf{x}_t \in \mathbb{R}^d, t \in [0, T]$  with the marginal probability  $p(\mathbf{x}_t)$ , where  $\bar{\mathbf{w}} \in \mathbb{R}^d$  is the reverse Wiener process.

508 **Proof E.1** For the clarity of the notation, we denote  $p(\mathbf{x}_t)$  as  $p(\mathbf{x}, t)$ , follow the Fokker-Planck  
509 equation [49], we can convert the SDE  $d\mathbf{x}_t = f(t)\mathbf{x}_t dt + g(t)d\bar{\mathbf{w}}$  to a partial differential equation  
510 Eq. (28) and Eq. (29).

$$\begin{aligned}
\frac{\partial}{\partial t}p(\mathbf{x}, t) &= -\sum_{i=1}^d \frac{\partial}{\partial \mathbf{x}_i} (\mathbf{f}_i(\mathbf{x}, t) \cdot p(\mathbf{x}, t)) + \frac{1}{2} \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} (g^2(t) \cdot p(\mathbf{x}, t)) \\
&= -\sum_{i=1}^d \frac{\partial}{\partial \mathbf{x}_i} (\mathbf{f}(\mathbf{x}, t) \cdot p(\mathbf{x}, t)) + \frac{\zeta}{2} \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} (g^2(t) \cdot p(\mathbf{x}, t)) \\
&\quad + \frac{1-\zeta}{2} \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} (g^2(t) \cdot p(\mathbf{x}, t)) \tag{28} \\
&= -\sum_{i=1}^d \frac{\partial}{\partial \mathbf{x}_i} (\mathbf{f}(\mathbf{x}, t) \cdot p(\mathbf{x}, t)) + \frac{\zeta}{2} \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} (g^2(t) \cdot p(\mathbf{x}, t)) \\
&\quad + \frac{1-\zeta}{2} \sum_{i=1}^d \frac{\partial}{\partial \mathbf{x}_i} (g^2(t) \cdot \frac{p(\mathbf{x}, t)}{p(\mathbf{x}, t)}) \nabla_{\mathbf{x}} p(\mathbf{x}, t)
\end{aligned}$$

511 To simplify the second-order partial derivative, in the Eq. (29), we apply the log-derivative trick:

$$\begin{aligned}
512 \log p(\mathbf{x}, t) \nabla_{\mathbf{x}} p(\mathbf{x}, t) &= \frac{\nabla_{\mathbf{x}} p(\mathbf{x}, t)}{p(\mathbf{x}, t)} \\
&= -\sum_{i=1}^d \frac{\partial}{\partial \mathbf{x}_i} (\mathbf{f}(\mathbf{x}, t) \cdot p(\mathbf{x}, t)) + \frac{\zeta}{2} \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} (g^2(t) \cdot p(\mathbf{x}, t)) \\
&\quad + \frac{1-\zeta}{2} \sum_{i=1}^d \frac{\partial}{\partial \mathbf{x}_i} ((g^2(t) \cdot \nabla_{\mathbf{x}} \log p(\mathbf{x}, t)) \cdot p(\mathbf{x}, t)) \tag{29} \\
&= -\sum_{i=1}^d \frac{\partial}{\partial \mathbf{x}_i} ((\mathbf{f}(\mathbf{x}, t) - \frac{1-\zeta}{2}g^2(t) \cdot \nabla_{\mathbf{x}} \log p(\mathbf{x}, t)) \cdot p(\mathbf{x}, t)) \\
&\quad + \frac{\zeta}{2} \sum_{i=1}^d \sum_{j=1}^d \frac{\partial^2}{\partial \mathbf{x}_i \partial \mathbf{x}_j} (g^2(t) \cdot p(\mathbf{x}, t))
\end{aligned}$$

513 Thus, we can convert the above results back to an SDE with the Fokker-Planck equation with  
514 randomness indicator  $\zeta$  in Eq. (30). We can see it will reduce to an ODE while  $\zeta = 0$  and SDE while  
515  $\zeta = 1$ .

$$d\mathbf{x}_t = [\mathbf{f}(\mathbf{x}_t, t) - \frac{1-\zeta}{2}g^2(t)\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)]dt + \sqrt{\zeta}\mathbf{g}(t)d\bar{\mathbf{w}} \tag{30}$$

516 ■

### 517 E.3.1 Backdoor Reversed SDE with Arbitrary Stochasticity

518 Since  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t) \approx q(\mathbf{x}_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)$ , we can replace  $\mathbf{x}_0$  of Eq. (25) with reparametrization  $\mathbf{x}_0 =$   
519  $\frac{\mathbf{x}'_t - \hat{\rho}(t)\mathbf{r} - \hat{\beta}(t)\epsilon_t}{\hat{\alpha}(t)}$  from Eq. (25). Note that since the marginal distribution  $q(\mathbf{x}'_t)$  follows Gaussian  
520 distribution, we replace the  $\epsilon_t$  with the normalized conditional score function  $-\hat{\beta}(t)\nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t|\mathbf{x}'_0)$   
521 as a kind of reparametrization trick.

$$\begin{aligned}
\mathbf{x}'_{t-1} &= a(t)\mathbf{x}'_t + b(t) \frac{\mathbf{x}'_t - \hat{\rho}(t)\mathbf{r} - \hat{\beta}(t)(-\hat{\beta}(t)\nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t|\mathbf{x}'_0))}{\hat{\alpha}(t)} + c(t)\mathbf{r} + s(t)\epsilon_t, \epsilon_t \sim \mathcal{N}(0, \mathbf{I}) \\
&= (a(t) + \frac{b(t)}{\hat{\alpha}(t)})\mathbf{x}'_t + (c(t) - \frac{b(t)\hat{\rho}(t)}{\hat{\alpha}(t)})\mathbf{r} - \frac{b(t)\hat{\beta}(t)}{\hat{\alpha}(t)}(-\hat{\beta}(t)\nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t|\mathbf{x}'_0)) + s(t)\epsilon_t \tag{31}
\end{aligned}$$

522 Then, based on Eq. (31), we approximate the dynamic  $d\mathbf{x}'_t$  with Taylor expansion as Eq. (32)

$$d\mathbf{x}'_t = \left[ \left( a(t) + \frac{b(t)}{\hat{\alpha}(t)} - 1 \right) \mathbf{x}'_t + \left( c(t) - \frac{b(t)\hat{\rho}(t)}{\hat{\alpha}(t)} \right) \mathbf{r} - \frac{b(t)\hat{\beta}(t)}{\hat{\alpha}(t)} \left( -\hat{\beta}(t) \nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t | \mathbf{x}'_0) \right) \right] dt + s(t) d\bar{\mathbf{w}} \quad (32)$$

523 With proper reorganization, we can express the SDE Eq. (32) as Eq. (33)

$$d\mathbf{x}'_t = \left[ F(t) \mathbf{x}'_t - G^2(t) \left( -\hat{\beta}(t) \nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t | \mathbf{x}'_0) - \frac{H(t)}{G^2(t)} \mathbf{r} \right) \right] dt + s(t) d\bar{\mathbf{w}} \quad (33)$$

524 We denote  $F(t) = a(t) + \frac{b(t)}{\hat{\alpha}(t)} - 1$ ,  $H(t) = c(t) - \frac{b(t)\hat{\rho}(t)}{\hat{\alpha}(t)}$ , and  $G(t) = \sqrt{\frac{b(t)\hat{\beta}(t)}{\hat{\alpha}(t)}}$ . Since we also  
 525 assume the forward process  $q(\mathbf{x}_t | \mathbf{x}_0)$  and  $q(\mathbf{x}'_t | \mathbf{x}'_0)$  are diffusion processes, thus the coefficient  $s(t)$   
 526 can be derived as  $s(t) = \sqrt{\hat{\beta}(t)G(t)} = \sqrt{\frac{b(t)}{\hat{\alpha}(t)}\hat{\beta}(t)}$ . Then, considering different stochasticity of  
 527 various samplers, we can apply Lemma 1 and introduce an additional stochasticity indicator  $\zeta \in [0, 1]$   
 528 in Eq. (34).

$$\begin{aligned} d\mathbf{x}'_t &= \left[ F(t) \mathbf{x}'_t - G^2(t) \left( -\hat{\beta}(t) \nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t | \mathbf{x}'_0) - \frac{H(t)}{G^2(t)} \mathbf{r} \right) \right] dt + s(t) d\bar{\mathbf{w}} \\ &= \left[ F(t) \mathbf{x}'_t - G^2(t) \left( -\hat{\beta}(t) \nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t | \mathbf{x}'_0) - \frac{H(t)}{G^2(t)} \mathbf{r} \right) - \frac{1-\zeta}{2} s^2(t) \nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t | \mathbf{x}'_0) \right] dt + \sqrt{\zeta} s(t) d\bar{\mathbf{w}} \\ &= \left[ F(t) \mathbf{x}'_t - \frac{1+\zeta}{2} G^2(t) \underbrace{\left( -\hat{\beta}(t) \nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t | \mathbf{x}'_0) - \frac{2H(t)}{(1+\zeta)G^2(t)} \mathbf{r} \right)}_{\text{Backdoor Score Function}} \right] dt + G(t) \sqrt{\zeta \hat{\beta}(t)} d\bar{\mathbf{w}} \end{aligned} \quad (34)$$

### 529 E.3.2 Learned Reversed SDE with Arbitrary Stochasticity

530 Since  $q(\mathbf{x}_{t-1} | \mathbf{x}_t) \approx q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ , we can replace  $\mathbf{x}_0$  of Eq. (25) with  $\mathbf{x}_0 = \frac{\mathbf{x}_t - \hat{\beta}(t)\epsilon_\theta(\mathbf{x}_t, t)}{\hat{\alpha}(t)}$ ,  
 531 which is derived from the reparametrization of the forward process  $\mathbf{x}_t = \hat{\alpha}(t)\mathbf{x}_0 + \hat{\beta}(t)\epsilon_t$  with the  
 532 replacement  $\epsilon_t$  with  $\epsilon_\theta(\mathbf{x}_t, t)$ .

$$\begin{aligned} \mathbf{x}_{t-1} &= a(t)\mathbf{x}_t + b(t) \frac{\mathbf{x}_t - \hat{\beta}(t)\epsilon_\theta(\mathbf{x}_t, t)}{\hat{\alpha}(t)} + s(t)\epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, \mathbf{I}) \\ &= \left( a(t) + \frac{b(t)}{\hat{\alpha}(t)} \right) \mathbf{x}_t - \frac{b(t)\hat{\beta}(t)}{\hat{\alpha}(t)} \epsilon_\theta(\mathbf{x}_t, t) + s(t)\epsilon_t \end{aligned} \quad (35)$$

533 Then, according to Eq. (35), we approximate the dynamic  $d\mathbf{x}_t$  with Taylor expansion as Eq. (36)

$$d\mathbf{x}_t = \left[ \left( a(t) + \frac{b(t)}{\hat{\alpha}(t)} - 1 \right) \mathbf{x}_t - \frac{b(t)\hat{\beta}(t)}{\hat{\alpha}(t)} \epsilon_\theta(\mathbf{x}_t, t) \right] dt + s(t) d\bar{\mathbf{w}} \quad (36)$$

534 With proper reorganization, we can express the SDE Eq. (36) with  $F(t) = a(t) + \frac{b(t)}{\hat{\alpha}(t)} - 1$ ,  
 535  $G(t) = \sqrt{\frac{b(t)\hat{\beta}(t)}{\hat{\alpha}(t)}}$ , and  $s(t) = \sqrt{\frac{b(t)}{\hat{\alpha}(t)}\hat{\beta}(t)}$  as Eq. (37).

$$d\mathbf{x}_t = \left[ F(t) \mathbf{x}_t - G^2(t) \epsilon_\theta(\mathbf{x}_t, t) \right] dt + s(t) d\bar{\mathbf{w}} \quad (37)$$

536 Then, we also consider arbitrary stochasticity and introduce an additional stochasticity indicator  
 537  $\zeta \in [0, 1]$  with Lemma 1. As we use a diffusion model  $\epsilon_\theta$  as an approximation for the normalized  
 538 score function:  $\epsilon_\theta(\mathbf{x}_t, t) = -\hat{\beta}(t) \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$ , we can derive the learned reversed SDE with  
 539 arbitrary stochasticity in Eq. (38).

$$d\mathbf{x}_t = \left[ F(t) \mathbf{x}_t - \frac{1+\zeta}{2} G^2(t) \epsilon_\theta(\mathbf{x}_t, t) \right] dt + G(t) \sqrt{\zeta \hat{\beta}(t)} d\bar{\mathbf{w}} \quad (38)$$

### 540 E.3.3 Loss Function of the Backdoor Diffusion Models

541 Based on the above results, we can formulate a score-matching problem based on Eq. (34) and  
 542 Eq. (38) as Eq. (39). The loss function Eq. (39) is also known as denoising-score-matching loss



543 [47], which is a surrogate of the score-matching problem since the score function  $\nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t)$  is  
 544 intractable.

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}'_t, \mathbf{x}'_0} \left[ \left\| \left( -\hat{\beta}(t) \nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t | \mathbf{x}'_0) - \frac{2H(t)}{(1+\zeta)G^2(t)} \mathbf{r} \right) - \epsilon_\theta(\mathbf{x}'_t, t) \right\|^2 \right] \\ & \propto \left\| \epsilon - \frac{2H(t)}{(1+\zeta)G^2(t)} \mathbf{r}(\mathbf{x}_0, \mathbf{g}) - \epsilon_\theta(\mathbf{x}'_t(\mathbf{x}'_0, \mathbf{r}(\mathbf{x}_0, \mathbf{g}), \epsilon), t) \right\|^2 \end{aligned} \quad (39)$$

545 Thus, we can finally write down the backdoor loss function Eq. (40).

$$L_p(\mathbf{x}, t, \epsilon, \mathbf{g}, \mathbf{y}, \zeta) := \left\| \epsilon - \frac{2H(t)}{(1+\zeta)G^2(t)} \mathbf{r}(\mathbf{x}, \mathbf{g}) - \epsilon_\theta(\mathbf{x}'_t(\mathbf{y}, \mathbf{r}(\mathbf{x}, \mathbf{g}), \epsilon), t) \right\|^2 \quad (40)$$

## 546 E.4 The Derivation of Conditional Diffusion Models

547 We will expand our framework to conditional generation in this section. In Appendix E.4.1, we  
 548 will start with the negative-log likelihood (NLL) and derive the variational lower bound (VLBO).  
 549 Next, in Appendix E.4.2, we decompose the VLBO into three components and focus on the most  
 550 important one. In Appendix E.4.3, based on previous sections, we will derive the clean loss function  
 551 for the conditional diffusion models. The last section Appendix E.4.4 will combine the results of  
 552 Appendix E.1 and Appendix E.2 and derive the backdoor loss functions for the conditional diffusion  
 553 models and various samplers.

### 554 E.4.1 Conditional Negative Log Likelihood (NLL)

555 To train a conditional diffusion model  $\epsilon_\theta(\mathbf{x}_0, \mathbf{c})$ , we will optimize the joint probability learned by  
 556 the model  $\arg \min_\theta -\log p_\theta(\mathbf{x}_0, \mathbf{c})$ . We denote  $\mathbf{c}$  as the condition, which can be prompt embedding  
 557 for the text-to-image generation, and the  $D_{\mathbf{x}_{i:T}}$  is the domain of random vectors  $\mathbf{x}_i, \dots, \mathbf{x}_T$ ,  $\mathbf{x}_t \in$   
 558  $\mathbb{R}^d$ ,  $t \in [i, T]$ ,  $i \leq T$ . Therefore, we can derive the conditional variational lower bound  $L_{VLB}^C$  as  
 559 Eq. (41).

$$\begin{aligned} -\log p_\theta(\mathbf{x}_0, \mathbf{c}) &= -\mathbb{E}_{q(\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0, \mathbf{c})] \\ &= -\mathbb{E}_{q(\mathbf{x}_0)} \left[ \log \int_{D_{\mathbf{x}_{1:T}}} p_\theta(\mathbf{x}_0, \dots, \mathbf{x}_T, \mathbf{c}) d\mathbf{x}_1 \dots d\mathbf{x}_T \right] \\ &= -\mathbb{E}_{q(\mathbf{x}_0)} \left[ \log \int_{\mathbf{x}_1 \dots \mathbf{x}_T} q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0) \frac{p_\theta(\mathbf{x}_0, \dots, \mathbf{x}_T, \mathbf{c})}{q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0)} d\mathbf{x}_1 \dots d\mathbf{x}_T \right] \\ &= -\mathbb{E}_{q(\mathbf{x}_0)} \left[ \log \mathbb{E}_{q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0)} \left[ \frac{p_\theta(\mathbf{x}_0, \dots, \mathbf{x}_T, \mathbf{c})}{q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0)} \right] \right] \\ &\leq -\mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)} \left[ \log \frac{p_\theta(\mathbf{x}_0, \dots, \mathbf{x}_T, \mathbf{c})}{q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0)} \right] = L_{VLB}^C \end{aligned} \quad (41)$$

### 560 E.4.2 Conditional Variational Lower Bound (VLBO)

561 In this section, we will further decompose the VLBO Eq. (41) and show that minimizing the  
 562 KL-divergence  $D_{KL}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c}))$  is our main objective in Eq. (43). For the  
 563 simplicity, we denote  $\mathbb{E}_{q(\mathbf{x}_0, \dots, \mathbf{x}_T)}$  as  $\mathbb{E}_q$ . With Markovian assumption, the latent  $\mathbf{x}_t$  at the timestep  $t$   
 564 only depends on the previous latent  $\mathbf{x}_{t-1}$  and the condition  $\mathbf{c}$ .

$$\begin{aligned} L_{VLB}^C &= -\mathbb{E}_q \left[ \log \frac{p_\theta(\mathbf{x}_0, \dots, \mathbf{x}_T, \mathbf{c})}{q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0)} \right] = \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0)}{p_\theta(\mathbf{x}_0, \dots, \mathbf{x}_T, \mathbf{c})} \right] \\ &= \mathbb{E}_q \left[ \log \frac{\prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_T, \mathbf{c}) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})} \right] \\ &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T, \mathbf{c}) + \sum_{t=1}^T \log \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})} \right] \\ &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T, \mathbf{c}) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})} + \log \frac{q(\mathbf{x}_1 | \mathbf{x}_0)}{p_\theta(\mathbf{x}_0 | \mathbf{x}_1, \mathbf{c})} \right] \\ &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T, \mathbf{c}) + \sum_{t=2}^T \log \left( \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})} \cdot \frac{q(\mathbf{x}_t | \mathbf{x}_0)}{q(\mathbf{x}_{t-1} | \mathbf{x}_0)} \right) + \log \frac{q(\mathbf{x}_1 | \mathbf{x}_0)}{p_\theta(\mathbf{x}_0 | \mathbf{x}_1, \mathbf{c})} \right] \end{aligned} \quad (42)$$

565

$$\begin{aligned}
&= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T, \mathbf{c}) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1, \mathbf{c})} \right] \\
&= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T, \mathbf{c}) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})} + \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{q(\mathbf{x}_1|\mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1|\mathbf{x}_0)}{p_\theta(\mathbf{x}_0|\mathbf{x}_1, \mathbf{c})} \right] \\
&= \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_T|\mathbf{x}_0)}{p_\theta(\mathbf{x}_T, \mathbf{c})} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})} - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1, \mathbf{c}) \right] \\
&= \mathbb{E}_q \left[ D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0)||p_\theta(\mathbf{x}_T, \mathbf{c})) + \sum_{t=2}^T D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})) - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1, \mathbf{c}) \right] \\
&= \mathbb{E}_q \left[ \mathcal{L}_T^C(\mathbf{x}_T, \mathbf{x}_0, \mathbf{c}) + \sum_{t=2}^T \mathcal{L}_t^C(\mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{x}_0, \mathbf{c}) - \mathcal{L}_0^C(\mathbf{x}_1, \mathbf{x}_0, \mathbf{c}) \right]
\end{aligned} \tag{43}$$

### 566 E.4.3 Clean Loss Function for the Conditional Diffusion Models

567 We define the learned reversed transition  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})$  as Eq. (44).

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c}) := \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, \mathbf{x}_0, t, \mathbf{c}), s^2(t)\mathbf{I}) \tag{44}$$

568 We plug in a conditional diffusion model  $\epsilon_\theta(\mathbf{x}_t, t, \mathbf{c})$  to replace the unconditional diffusion model  
569  $\epsilon_\theta(\mathbf{x}_t, t)$ .

$$\begin{aligned}
\mu_\theta(\mathbf{x}_t, \mathbf{x}_0, t, \mathbf{c}) &= \frac{k_t \hat{\beta}^2(t-1)}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \mathbf{x}_t + \frac{\hat{\alpha}(t-1)w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \left( \frac{1}{\hat{\alpha}(t)} (\mathbf{x}_t - \hat{\beta}(t)\epsilon_\theta(\mathbf{x}_t, t, \mathbf{c})) \right) \\
&= \frac{k_t \hat{\beta}^2(t-1)\hat{\alpha}(t) + \hat{\alpha}(t-1)w_t^2}{\hat{\alpha}(t)(k_t^2 \hat{\beta}^2(t-1) + w_t^2)} \mathbf{x}_t - \frac{\hat{\alpha}(t-1)w_t^2}{k_t^2 \hat{\beta}^2(t-1) + w_t^2} \frac{\hat{\beta}(t)}{\hat{\alpha}(t)} \epsilon_\theta(\mathbf{x}_t, t, \mathbf{c})
\end{aligned} \tag{45}$$

570 As a result, we use mean-matching as an approximation of the KL-divergence loss with Eq. (46).

$$D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})) \propto \|\mu_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, \mathbf{x}_0, t, \mathbf{c})\|^2 \propto \|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t, \mathbf{c})\|^2 \tag{46}$$

571 Finally, we can reorganize the Eq. (46) as Eq. (47), which is the clean loss function for the conditional  
572 diffusion models.

$$\mathcal{L}_c^C(\mathbf{x}, t, \epsilon, \mathbf{c}) := \|\epsilon - \epsilon_\theta(\mathbf{x}_t(\mathbf{x}, \epsilon), t, \mathbf{c})\|^2 \tag{47}$$

### 573 E.4.4 Loss Function of the Backdoor Conditional Diffusion Models

574 Based on the above results, we can further derive the learned conditional reversed SDE Eq. (48),  
575 while the backdoor one remains the same as Eq. (34), which is caused by the identical backdoor  
576 reversed transition  $q(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)$  of the KL-divergence loss.

$$d\mathbf{x}_t = \left[ F(t)\mathbf{x}_t - \frac{1+\zeta}{2}G^2(t)\epsilon_\theta(\mathbf{x}_t, t, \mathbf{c}) \right] dt + G(t)\sqrt{\zeta\hat{\beta}(t)}d\bar{\mathbf{w}} \tag{48}$$

577 According to the above results, we can formulate an image-trigger backdoor loss function based  
578 on Eq. (34) and Eq. (48) as Eq. (49). The loss function Eq. (49) is also known as denoising-score-  
579 matching loss [47], which is a surrogate of the score-matching problem since the score function  
580  $\nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t)$  is intractable. Here we denote the reparametrization  $\mathbf{x}'_t(\mathbf{x}, \mathbf{r}, \epsilon) = \hat{\alpha}(t)\mathbf{x} + \hat{\rho}(t)\mathbf{r} +$   
581  $\hat{\beta}(t)\epsilon$ .

$$\begin{aligned}
&\mathbb{E}_{\mathbf{x}'_t, \mathbf{x}'_0} \left[ \left\| (-\hat{\beta}(t)\nabla_{\mathbf{x}'_t} \log q(\mathbf{x}'_t|\mathbf{x}'_0) - \frac{2H(t)}{(1+\zeta)G^2(t)}\mathbf{r}) - \epsilon_\theta(\mathbf{x}_t, t, \mathbf{c}) \right\|^2 \right] \\
&\propto \left\| \epsilon - \frac{2H(t)}{(1+\zeta)G^2(t)}\mathbf{r}(\mathbf{x}_0, \mathbf{g}) - \epsilon_\theta(\mathbf{x}'_t(\mathbf{x}'_0, \mathbf{r}(\mathbf{x}_0, \mathbf{g}), \epsilon), t, \mathbf{c}) \right\|^2
\end{aligned} \tag{49}$$

582 Thus, we can finally write down the image-as-trigger backdoor loss function Eq. (50) for the  
583 conditional diffusion models.

$$\mathcal{L}_p^{CI}(\mathbf{x}, t, \epsilon, \mathbf{g}, \mathbf{y}, \mathbf{c}, \zeta) := \left\| \epsilon - \frac{2H(t)}{(1+\zeta)G^2(t)}\mathbf{r}(\mathbf{y}, \mathbf{g}) - \epsilon_\theta(\mathbf{x}'_t(\mathbf{y}, \mathbf{r}(\mathbf{x}, \mathbf{g}), \epsilon), t, \mathbf{c}) \right\|^2 \tag{50}$$

## 584 F Additional Experiments

### 585 F.1 Backdoor Attacks on DDPM with CIFAR10 Dataset

586 We will present experimental results for more backdoor trigger-target pairs and samplers, including  
 587 the LMSD sampler, which is implemented by the authors of EDM [20], in Fig. 7. The results of the ANCESTRAL sampler come from [7].

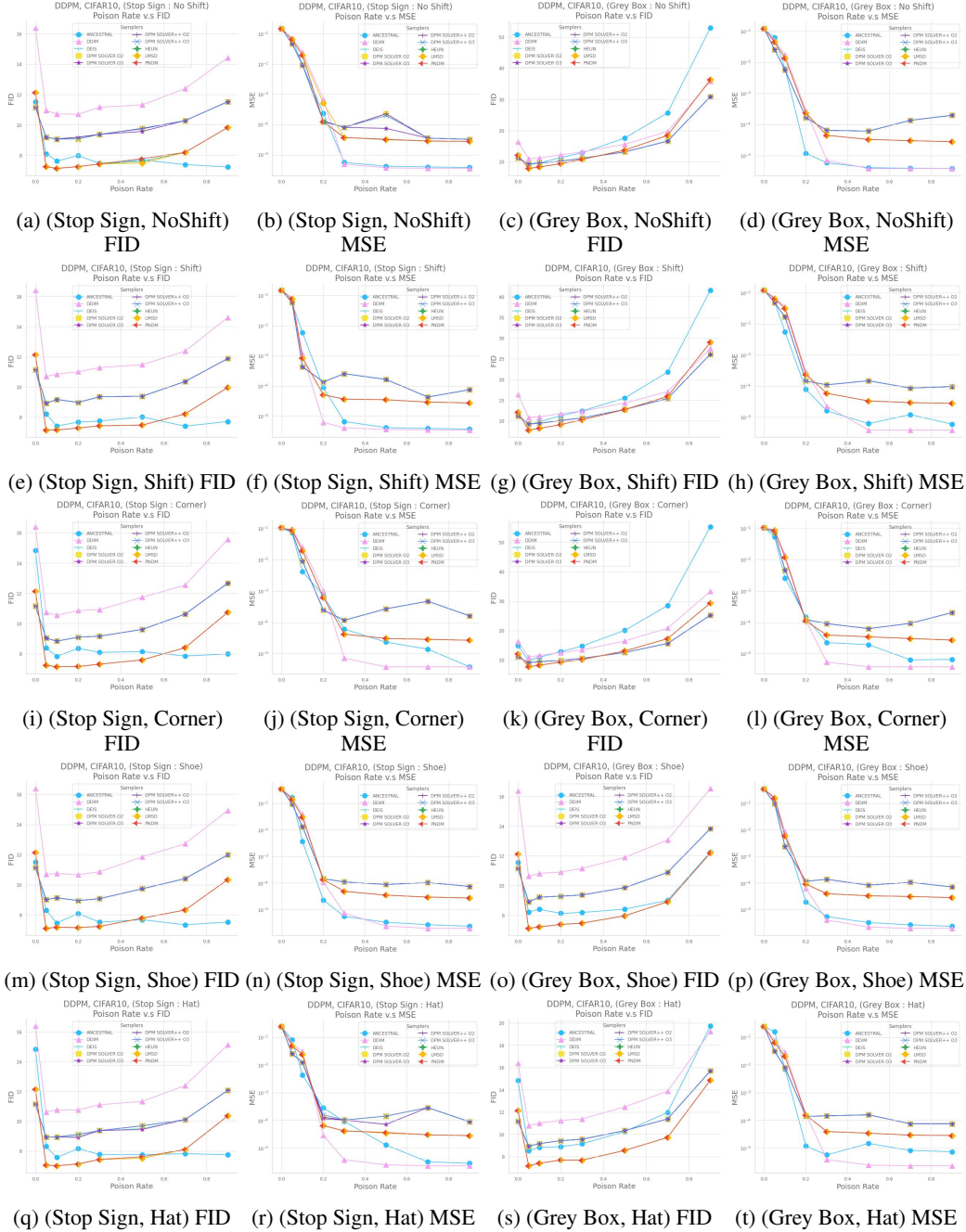
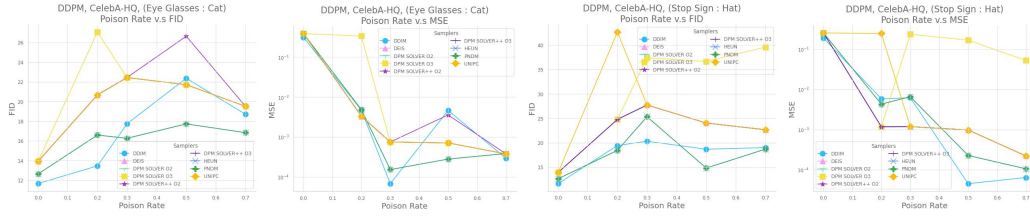


Figure 7: FID and MSE scores of various samplers and poison rates for DDPM [11] and the CIFAR10 dataset. We express trigger-target pairs as (trigger, target).

589 **F.2 Backdoor Attacks on DDPM with CelebA-HQ Dataset**

590 We evaluate our method with more samplers and backdoor trigger-target pairs: (Stop Sign, Hat) and (Eyeglasses, Cat) in Fig. 8. Note that the results of the ANCESTRAL sampler come from [7].



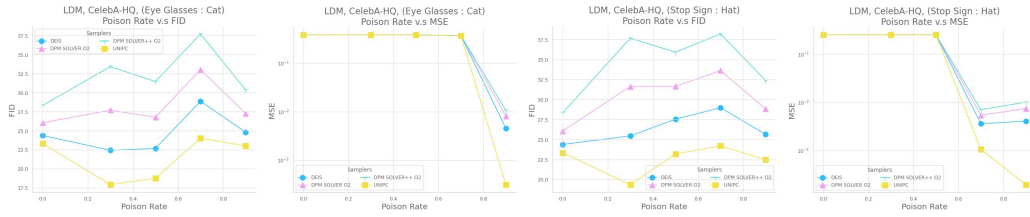
(a) (Eyeglasses, Cat) FID (b) (Eyeglasses, Cat) MSE (c) (Stop Sign, Hat) FID (d) (Stop Sign, Hat) MSE

Figure 8: FID and MSE scores of various samplers and poison rates for the DDPM [11] and the CelebA-HQ dataset. We express trigger-target pairs as (trigger, target).

591

592 **F.3 Backdoor Attacks on Latent Diffusion Models (LDM)**

593 We also attack the latent diffusion model [39] downloaded from Huggingface (*CompVis/lDM-celebahq-*  
 594 *256*), which is pre-trained on CelebA-HQ [27]. The pre-trained latent diffusion models (LDM) [39]  
 595 are trained on CelebA-HQ with  $512 \times 512$  resolution and  $64 \times 64$  latent space. We fine-tune them  
 596 with learning rate  $2e-4$  and batch size 16 for 2000 epochs. We examine our method with trigger-target  
 597 pair: (Eyeglasses, Cat) and (Stop Sign, Hat) and illustrate the FID and MSE score in Fig. 9. As the  
 598 Fig. 9 shows, the LDM can be backdoored successfully for the trigger-target pairs: (Stop Sign, Hat)  
 599 with 70% poison rate. Meanwhile, for the trigger-target pair: (Eye Glasses, Cat) and 90% poison  
 600 rate, the FID scores only slightly increase by about 7.2% at most. As for the trigger-target pair: (Stop  
 601 Sign, Hat), although the FID raises higher than (Eye Glasses, Cat), we believe longer training can  
 enhance their utility.



(a) (Eyeglasses, Cat) FID (b) (Eyeglasses, Cat) MSE (c) (Stop Sign, Hat) FID (d) (Stop Sign, Hat) MSE

Figure 9: FID and MSE scores of various samplers and poison rates for the latent diffusion model (LDM) [39] and the CelebA-HQ dataset. We express trigger-target pairs as (trigger, target).

602

603 **F.4 Backdoor Attacks on Score-Based Models**

604 We trained the score-based model: NCSN [49, 47, 48] on the CIFAR10 dataset with the same model  
 605 architecture as the DDPM [11] by ourselves for 800 epochs and set the learning rate as  $1e-4$  and  
 606 batch size as 128. The FID score of the clean model generated by predictor-corrector samplers  
 607 (SCORE-SDE-VE [49]) for the variance explode models [49] is about 10.87. For the backdoor, we  
 608 fine-tune the pre-trained model with the learning rate  $2e-5$  and batch size 128 for 46875 steps. To  
 609 enhance the backdoor specificity and utility, we augment Gaussian noise into the training dataset,  
 610 which means the poisoned image  $r$  will be replaced by a pure trigger  $g$ . The augmentation can let the  
 611 model learn to activate the backdoor even if there are no context images. We present our results in  
 612 Fig. 10 and can see with 70% augment rate, our method can achieve 70% attack success rate based on  
 613 Fig. 10c as the FID score increases by 12.7%. Note that the augment rate is computed by the number  
 614 of augmented Gaussian noises / the size of the original training dataset.

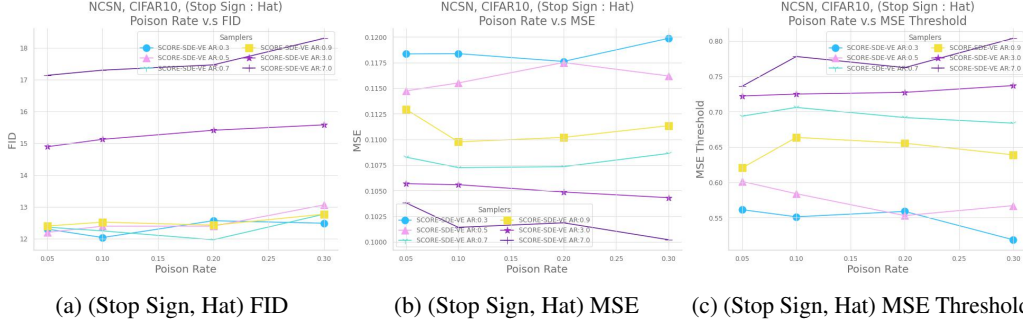


Figure 10: FID and MSE scores of various samplers and poison rates for the score-based model (NCSN) [47, 48, 49] and the CIFAR10 dataset. We express trigger-target pairs as (trigger, target). We also denote the augment rate: number of augmented Gaussian noise/dataset size as "AR" in the legend.

## 615 F.5 Inference-Time Clipping Defense

616 We evaluate the inference-time clipping defense on the CIFAR10 dataset with triggers: Grey Box  
 617 and Stop Sign and targets: NoShift, Shift, Corner, Shoe, and Hat in Fig. 11. The results of the  
 618 ANCESTRAL sampler are from [7]. We can see that inference-time clipping is still not effective for  
 619 most ODE samplers.

## 620 F.6 VillanDiffusion on the Inpaint Tasks

621 Similar to [7], we also evaluate our method on the inpainting tasks with various samplers. We design  
 622 3 kinds of different corruptions: **Blur**, **Line**, **Box**. **Blur** means we add Gaussian noise  $\mathcal{N}(0, 0.3)$  to  
 623 corrupt the images. **Line** and **Box** mean we crop part of the image and ask the diffusion models to  
 624 recover the missing area. We use VillanDiffusion trained on the trigger: Stop Sign and target: Shoe  
 625 and Hat with poison rate: 20%. During inpainting, we apply UniPC [54], DEIS [52], DPM Solver  
 626 [28], and DPM Solver++ [29] samplers with 50 sampling steps. To evaluate the recovery quality,  
 627 we generate 1024 images and use LPIPS [53] score to measure the similarity between the covered  
 628 images and ground-truth images. We illustrate our results in Fig. 12. We can see our method achieves  
 629 both high utility and high specificity.

## 630 G Numerical Results

### 631 G.1 Image-Trigger Backdoor Attacks on Unconditional Diffusion Models

#### 632 G.1.1 Backdoor Attacks on DDPM with CIFAR10 Dataset

633 We present the numerical results of the trigger: Stop Sign and targets: NoShift, Shift, Corner, Shoe,  
 634 and Hat in Table 2, Table 3, Table 4, Table 5, and Table 6 respectively. As for trigger Grey Box,  
 635 we also show the results for the targets: NoShift, Shift, Corner, Shoe, and Hat in Table 7, Table 8,  
 636 Table 9, Table 10, and Table 11.

#### 637 G.1.2 Backdoor Attacks on DDPM with CelebA-HQ Dataset

638 We show the numerical results for the trigger-target pairs: (Eyeglasses, Cat) and (Stop Sign, Hat) in  
 639 Table 12 and Table 13 respectively.

### 640 G.2 Backdoor Attacks on Latent Diffusion Models (LDM)

641 We show the experiment results of the trigger-target pair: (Eye Glasses, Cat) in Table 14 and (Stop  
 642 Sign, Hat) in Table 15.

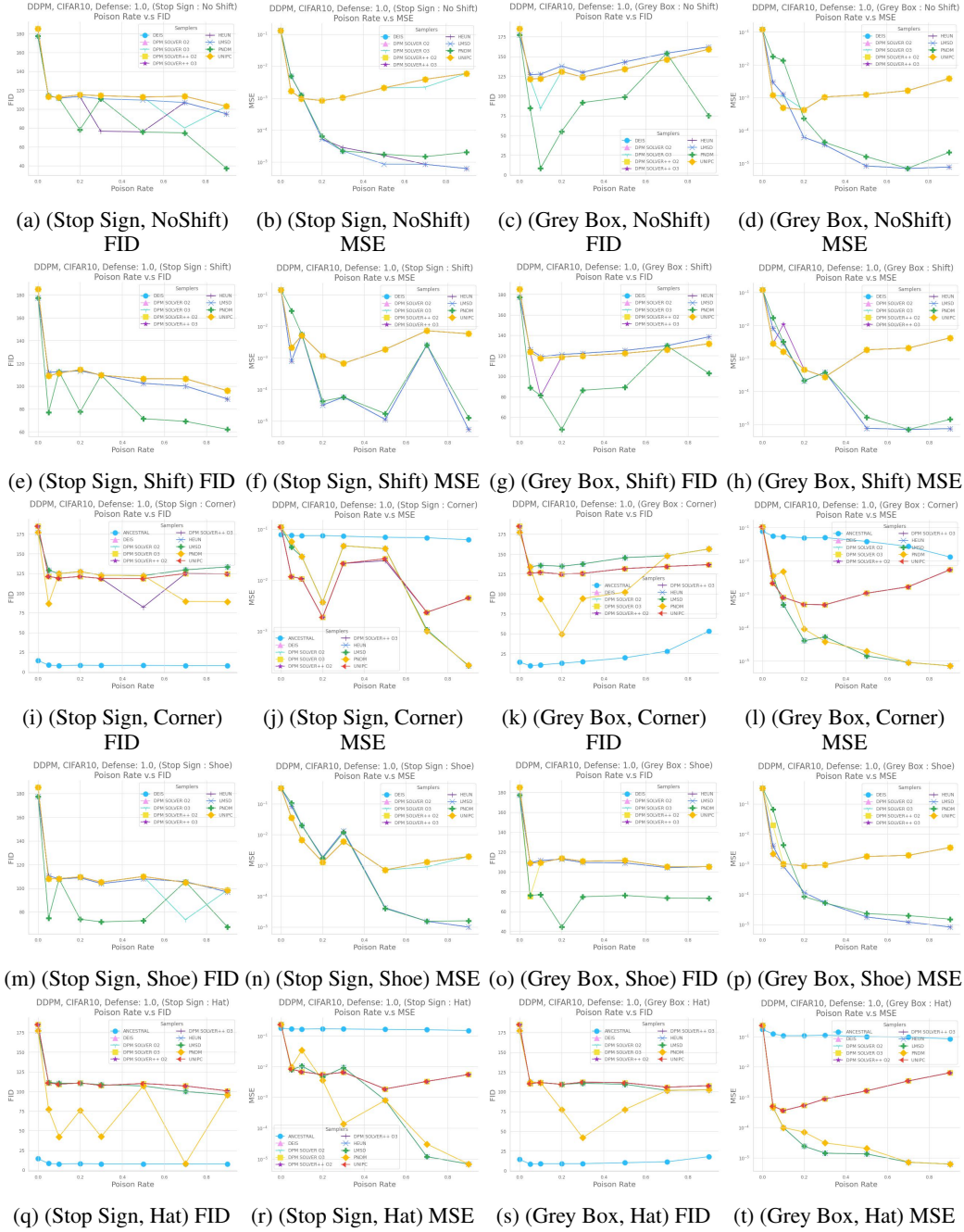
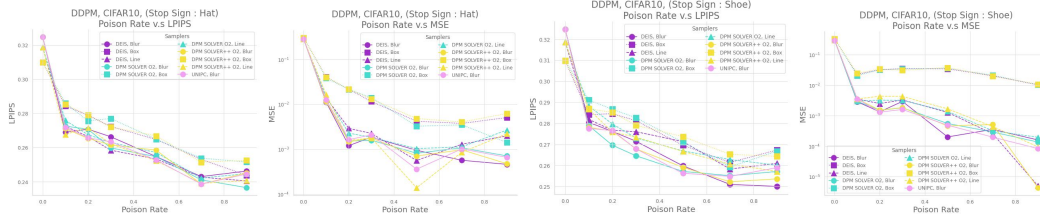


Figure 11: FID and MSE scores of various samplers and poison rates with inference-time defense [7]. We evaluate the defense on the DDPM and the CIFAR10 dataset and express trigger-target pairs as (trigger, target).



(a) (Stop Sign,Hat) LPIPS (b) (Stop Sign,Hat) MSE (c) (Stop Sign,Shoe) LPIPS (d) (Stop Sign,Shoe) MSE

Figure 12: LPIPS and MSE scores of various samplers and poison rates for the 3 kinds of inpainting tasks: **Blur**, **Line**, and **Box**. The backdoor model is DDPM trained on the CIFAR10 dataset. We express trigger-target pairs as (trigger, target).

Table 2: DDPM backdoor on CIFAR10 Dataset with Trigger: Stop Sign, target: No Shift

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	11.52	8.09	7.62	7.97	7.46	7.68	7.38	7.22
	MSE	1.48E-1	6.81E-2	9.47E-3	2.35E-4	5.59E-6	4.19E-6	3.96E-6	3.80E-6
	SSIM	6.84E-4	4.35E-1	9.18E-1	9.97E-1	9.99E-1	9.98E-1	9.98E-1	9.98E-1
UNIPC	FID	11.15	9.18	9.07	9.18	9.37	9.76	10.28	11.53
	MSE	1.48E-1	4.76E-2	9.37E-3	1.30E-4	8.05E-5	2.27E-4	3.56E-5	3.24E-5
	SSIM	8.13E-4	5.27E-1	8.57E-1	9.76E-1	9.81E-1	9.74E-1	9.86E-1	9.84E-1
DPM. O2	FID	11.15	9.18	9.07	9.07	9.37	9.76	10.28	11.53
	MSE	1.48E-1	4.76E-2	9.37E-3	1.24E-4	8.05E-5	2.27E-4	3.56E-5	3.24E-5
	SSIM	8.13E-4	5.27E-1	8.57E-1	9.80E-1	9.81E-1	9.74E-1	9.86E-1	9.84E-1
DPM. O3	FID	11.15	9.18	9.07	9.18	9.37	9.57	10.28	11.53
	MSE	1.48E-1	4.76E-2	9.37E-3	1.30E-4	8.05E-5	7.48E-5	3.56E-5	3.24E-5
	SSIM	8.13E-4	5.27E-1	8.57E-1	9.76E-1	9.81E-1	9.80E-1	9.86E-1	9.84E-1
DPM++. O2	FID	11.15	9.18	9.07	9.07	9.37	9.76	10.28	11.53
	MSE	1.48E-1	4.76E-2	9.37E-3	1.30E-4	8.05E-5	2.27E-4	3.56E-5	3.24E-5
	SSIM	8.13E-4	5.27E-1	8.57E-1	9.76E-1	9.81E-1	9.74E-1	9.86E-1	9.84E-1
DPM++. O3	FID	11.15	9.18	9.07	9.07	9.37	9.73	10.28	11.53
	MSE	1.48E-1	4.76E-2	9.37E-3	1.24E-4	8.05E-5	1.99E-4	3.56E-5	3.24E-5
	SSIM	8.13E-4	5.27E-1	8.57E-1	9.80E-1	9.81E-1	9.76E-1	9.86E-1	9.84E-1
DEIS	FID	11.15	9.18	9.07	9.07	9.37	9.57	10.28	11.53
	MSE	1.48E-1	4.76E-2	9.37E-3	1.28E-4	8.05E-5	7.48E-5	3.56E-5	3.24E-5
	SSIM	8.13E-4	5.27E-1	8.57E-1	9.76E-1	9.81E-1	9.80E-1	9.86E-1	9.84E-1
DDIM	FID	16.39	10.95	10.71	10.70	11.16	11.32	12.40	14.43
	MSE	1.48E-1	7.14E-2	2.47E-2	6.84E-4	4.95E-6	3.70E-6	3.58E-6	3.51E-6
	SSIM	8.92E-4	3.74E-1	7.63E-1	9.92E-1	1.00E+0	9.99E-1	9.99E-1	9.99E-1
PNDM	FID	12.14	7.24	7.13	7.25	7.42	7.77	8.18	9.83
	MSE	1.48E-1	6.55E-2	1.97E-2	1.23E-4	3.74E-5	3.25E-5	2.86E-5	2.76E-5
	SSIM	8.23E-4	4.11E-1	7.91E-1	9.81E-1	9.83E-1	9.83E-1	9.83E-1	9.83E-1
HEUN	FID	12.14	7.24	7.13	7.25	7.42	7.60	8.18	9.83
	MSE	1.48E-1	6.55E-2	1.97E-2	1.23E-4	3.74E-5	3.16E-5	2.86E-5	2.76E-5
	SSIM	8.23E-4	4.11E-1	7.91E-1	9.81E-1	9.83E-1	9.83E-1	9.83E-1	9.83E-1
LMSD	FID	12.14	7.24	7.13	7.24	7.42	7.47	8.18	9.83
	MSE	1.48E-1	6.55E-2	1.97E-2	5.02E-4	3.74E-5	3.23E-5	2.86E-5	2.76E-5
	SSIM	8.23E-4	4.11E-1	7.91E-1	9.76E-1	9.83E-1	9.83E-1	9.83E-1	9.83E-1

Table 3: DDPM backdoor on CIFAR10 Dataset with Trigger: Stop Sign, target: Shift

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	11.16	8.21	7.42	7.68	7.76	8.02	7.42	7.72
	MSE	1.48E-1	5.68E-2	5.91E-3	8.96E-5	6.73E-6	4.23E-6	3.96E-6	3.80E-6
	SSIM	4.24E-4	5.73E-1	9.56E-1	9.99E-1	9.99E-1	9.99E-1	9.99E-1	9.99E-1
UNIPC	FID	11.15	8.92	9.18	8.98	9.36	9.39	10.37	11.89
	MSE	1.48E-1	5.92E-2	4.37E-4	1.37E-4	2.58E-4	1.66E-4	4.35E-5	7.59E-5
	SSIM	4.25E-4	4.92E-1	9.70E-1	9.85E-1	9.81E-1	9.86E-1	9.88E-1	9.89E-1
DPM. O2	FID	11.15	8.92	9.18	8.98	9.36	9.39	10.37	11.89
	MSE	1.48E-1	5.92E-2	4.37E-4	1.37E-4	2.58E-4	1.66E-4	4.35E-5	7.59E-5
	SSIM	4.25E-4	4.92E-1	9.70E-1	9.85E-1	9.81E-1	9.86E-1	9.88E-1	9.89E-1
DPM. O3	FID	11.15	8.92	9.18	8.98	9.36	9.39	10.37	11.89
	MSE	1.48E-1	5.92E-2	4.37E-4	1.37E-4	2.58E-4	1.66E-4	4.35E-5	7.59E-5
	SSIM	4.25E-4	4.92E-1	9.70E-1	9.85E-1	9.81E-1	9.86E-1	9.88E-1	9.89E-1
DPM++. O2	FID	11.15	8.92	9.18	8.98	9.36	9.39	10.37	11.89
	MSE	1.48E-1	5.92E-2	4.37E-4	1.37E-4	2.58E-4	1.66E-4	4.35E-5	7.59E-5
	SSIM	4.25E-4	4.92E-1	9.70E-1	9.85E-1	9.81E-1	9.86E-1	9.88E-1	9.89E-1
DPM++. O3	FID	11.15	8.92	9.18	8.98	9.36	9.39	10.37	11.89
	MSE	1.48E-1	5.92E-2	4.37E-4	1.37E-4	2.58E-4	1.66E-4	4.35E-5	7.59E-5
	SSIM	4.25E-4	4.92E-1	9.70E-1	9.85E-1	9.81E-1	9.86E-1	9.88E-1	9.89E-1
DEIS	FID	11.15	8.92	9.18	8.98	9.36	9.39	10.37	11.89
	MSE	1.48E-1	5.92E-2	4.37E-4	1.37E-4	2.58E-4	1.66E-4	4.35E-5	7.59E-5
	SSIM	4.25E-4	4.92E-1	9.70E-1	9.85E-1	9.81E-1	9.86E-1	9.88E-1	9.89E-1
DDIM	FID	16.39	10.70	10.85	11.01	11.29	11.48	12.37	14.60
	MSE	1.48E-1	8.25E-2	1.18E-3	6.30E-6	4.15E-6	3.67E-6	3.54E-6	3.48E-6
	SSIM	4.32E-4	3.43E-1	9.84E-1	1.00E+0	1.00E+0	9.99E-1	9.99E-1	9.99E-1
PNDM	FID	12.14	7.15	7.16	7.29	7.44	7.48	8.23	9.97
	MSE	1.48E-1	7.69E-2	8.51E-4	5.18E-5	3.72E-5	3.53E-5	2.96E-5	2.77E-5
	SSIM	4.25E-4	3.78E-1	9.76E-1	9.89E-1	9.89E-1	9.90E-1	9.90E-1	9.90E-1
HEUN	FID	12.14	7.15	7.16	7.29	7.44	7.48	8.23	9.97
	MSE	1.48E-1	7.69E-2	8.51E-4	5.18E-5	3.72E-5	3.53E-5	2.96E-5	2.77E-5
	SSIM	4.25E-4	3.78E-1	9.76E-1	9.89E-1	9.89E-1	9.90E-1	9.90E-1	9.90E-1
LMSD	FID	12.14	7.15	7.16	7.29	7.44	7.48	8.23	9.97
	MSE	1.48E-1	7.69E-2	8.51E-4	5.18E-5	3.72E-5	3.53E-5	2.96E-5	2.77E-5
	SSIM	4.25E-4	3.78E-1	9.76E-1	9.89E-1	9.89E-1	9.90E-1	9.90E-1	9.90E-1

### 643 G.3 Backdoor Attacks on Score-Based Models

644 We provide the numerical results for the trigger-target pair: (Stop Sign, Hat) in the Table 16.

### 645 G.4 Caption-Trigger Backdoor Attacks on Text-to-Image DMs

646 We will present the numerical results of the Pokemon Caption dataset in Table 17 and Table 18. For  
647 the CelebA-HQ-Dialog dataset, we will show them in Table 19 and Table 20.

### 648 G.5 Inference-Time Clipping Defense

649 For the trigger Stop Sign, we present the numerical results of the inference-time clipping defense  
650 with targets: NoShift, Shift, Corner, Shoe, and Hat in Table 21, Table 22, Table 23, Table 24, and  
651 Table 25 respectively. As for the trigger Grey Box, we also show our results of the targets: NoShift,  
652 Shift, Corner, Shoe, and Hat in Table 26, Table 27, Table 28, Table 29, and Table 30 respectively.

### 653 G.6 VillanDiffusion on the Inpaint Tasks

654 For the trigger Stop Sign, we present the numerical results of the inpainting tasks: **Blur**, **Line**, and  
655 **Box** with targets Hat and Shoe in Table 32 and Table 31 respectively.



Table 4: DDPM backdoor on CIFAR10 Dataset with Trigger: Stop Sign, target: Corner

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	14.83	8.38	7.83	8.35	8.08	8.14	7.85	7.98
	MSE	1.06E-1	7.22E-2	4.20E-3	7.09E-4	6.13E-5	2.37E-5	1.41E-5	3.85E-6
	SSIM	9.85E-4	2.65E-1	9.49E-1	9.89E-1	9.97E-1	9.97E-1	9.97E-1	9.97E-1
UNIPC	FID	11.15	9.03	8.83	9.09	9.16	9.61	10.61	12.67
	MSE	1.06E-1	7.85E-2	9.05E-3	2.47E-4	1.17E-4	2.73E-4	4.75E-4	1.63E-4
	SSIM	1.09E-3	1.34E-1	7.52E-1	9.38E-1	9.56E-1	9.53E-1	9.57E-1	9.72E-1
DPM. O2	FID	11.15	9.03	8.83	9.09	9.16	9.61	10.61	12.67
	MSE	1.06E-1	7.85E-2	9.05E-3	2.47E-4	1.17E-4	2.73E-4	4.75E-4	1.63E-4
	SSIM	1.09E-3	1.34E-1	7.52E-1	9.38E-1	9.56E-1	9.53E-1	9.57E-1	9.72E-1
DPM. O3	FID	11.15	9.03	8.83	9.09	9.16	9.61	10.61	12.67
	MSE	1.06E-1	7.85E-2	9.05E-3	2.47E-4	1.17E-4	2.73E-4	4.75E-4	1.63E-4
	SSIM	1.09E-3	1.34E-1	7.52E-1	9.38E-1	9.56E-1	9.53E-1	9.57E-1	9.72E-1
DPM++. O2	FID	11.15	9.03	8.83	9.09	9.16	9.61	10.61	12.67
	MSE	1.06E-1	7.85E-2	9.05E-3	2.47E-4	1.17E-4	2.73E-4	4.75E-4	1.63E-4
	SSIM	1.09E-3	1.34E-1	7.52E-1	9.38E-1	9.56E-1	9.53E-1	9.57E-1	9.72E-1
DPM++. O3	FID	11.15	9.03	8.83	9.09	9.16	9.61	10.61	12.67
	MSE	1.06E-1	7.85E-2	9.05E-3	2.47E-4	1.17E-4	2.73E-4	4.75E-4	1.63E-4
	SSIM	1.09E-3	1.34E-1	7.52E-1	9.38E-1	9.56E-1	9.53E-1	9.57E-1	9.72E-1
DEIS	FID	11.15	9.03	8.83	9.09	9.16	9.61	10.61	12.67
	MSE	1.06E-1	7.85E-2	9.05E-3	2.47E-4	1.17E-4	2.73E-4	4.75E-4	1.63E-4
	SSIM	1.09E-3	1.34E-1	7.52E-1	9.38E-1	9.56E-1	9.53E-1	9.57E-1	9.72E-1
DDIM	FID	16.39	10.74	10.54	10.85	10.92	11.74	12.53	15.57
	MSE	1.06E-1	9.16E-2	2.54E-2	1.05E-3	7.27E-6	3.84E-6	3.84E-6	3.84E-6
	SSIM	1.12E-3	6.38E-2	6.36E-1	9.79E-1	1.00E+0	1.00E+0	9.99E-1	9.98E-1
PNDM	FID	12.14	7.22	7.14	7.15	7.31	7.59	8.39	10.74
	MSE	1.06E-1	8.87E-2	1.94E-2	6.28E-4	4.24E-5	3.09E-5	2.89E-5	2.70E-5
	SSIM	1.08E-3	7.93E-2	6.84E-1	9.58E-1	9.73E-1	9.76E-1	9.75E-1	9.76E-1
HEUN	FID	12.14	7.22	7.14	7.15	7.31	7.59	8.39	10.74
	MSE	1.06E-1	8.87E-2	1.94E-2	6.28E-4	4.24E-5	3.09E-5	2.89E-5	2.70E-5
	SSIM	1.08E-3	7.93E-2	6.84E-1	9.58E-1	9.73E-1	9.76E-1	9.75E-1	9.76E-1
LMSD	FID	12.14	7.22	7.14	7.15	7.31	7.59	8.39	10.74
	MSE	1.06E-1	8.87E-2	1.94E-2	6.28E-4	4.24E-5	3.09E-5	2.89E-5	2.70E-5
	SSIM	1.08E-3	7.93E-2	6.84E-1	9.58E-1	9.73E-1	9.76E-1	9.75E-1	9.76E-1

Table 5: DDPM backdoor on CIFAR10 Dataset with Trigger: Stop Sign, target: Shoe

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	11.52	8.33	7.47	8.10	7.52	7.69	7.35	7.54
	MSE	3.38E-1	1.66E-1	3.61E-3	2.30E-5	5.62E-6	3.35E-6	2.72E-6	2.39E-6
	SSIM	1.69E-4	4.20E-1	9.85E-1	9.99E-1	1.00E+0	1.00E+0	1.00E+0	1.00E+0
UNIPC	FID	11.15	9.03	9.14	8.96	9.09	9.74	10.41	12.00
	MSE	3.38E-1	9.33E-2	1.31E-2	1.47E-4	1.10E-4	8.83E-5	1.04E-4	7.49E-5
	SSIM	2.15E-4	5.91E-1	9.13E-1	9.89E-1	9.91E-1	9.92E-1	9.91E-1	9.92E-1
DPM. O2	FID	11.15	9.03	9.14	8.96	9.09	9.74	10.41	12.00
	MSE	3.38E-1	9.33E-2	1.31E-2	1.47E-4	1.10E-4	8.83E-5	1.04E-4	7.49E-5
	SSIM	2.15E-4	5.91E-1	9.13E-1	9.89E-1	9.91E-1	9.92E-1	9.91E-1	9.92E-1
DPM. O3	FID	11.15	9.03	9.14	8.96	9.09	9.74	10.41	12.00
	MSE	3.38E-1	9.33E-2	1.31E-2	1.47E-4	1.10E-4	8.83E-5	1.04E-4	7.49E-5
	SSIM	2.15E-4	5.91E-1	9.13E-1	9.89E-1	9.91E-1	9.92E-1	9.91E-1	9.92E-1
DPM++. O2	FID	11.15	9.03	9.14	8.96	9.09	9.74	10.41	12.00
	MSE	3.38E-1	9.33E-2	1.31E-2	1.47E-4	1.10E-4	8.83E-5	1.04E-4	7.49E-5
	SSIM	2.15E-4	5.91E-1	9.13E-1	9.89E-1	9.91E-1	9.92E-1	9.91E-1	9.92E-1
DPM++. O3	FID	11.15	9.03	9.14	8.96	9.09	9.74	10.41	12.00
	MSE	3.38E-1	9.33E-2	1.31E-2	1.47E-4	1.10E-4	8.83E-5	1.04E-4	7.49E-5
	SSIM	2.15E-4	5.91E-1	9.13E-1	9.89E-1	9.91E-1	9.92E-1	9.91E-1	9.92E-1
DEIS	FID	11.15	9.03	9.14	8.96	9.09	9.74	10.41	12.00
	MSE	3.38E-1	9.33E-2	1.31E-2	1.47E-4	1.10E-4	8.83E-5	1.04E-4	7.49E-5
	SSIM	2.15E-4	5.91E-1	9.13E-1	9.89E-1	9.91E-1	9.92E-1	9.91E-1	9.92E-1
DDIM	FID	16.39	10.71	10.75	10.68	10.87	11.86	12.73	14.94
	MSE	3.37E-1	1.56E-1	3.96E-2	1.09E-4	7.39E-6	2.42E-6	2.00E-6	1.98E-6
	SSIM	2.40E-4	3.97E-1	8.14E-1	9.99E-1	1.00E+0	1.00E+0	1.00E+0	1.00E+0
PNDM	FID	12.14	7.12	7.20	7.17	7.25	7.79	8.33	10.35
	MSE	3.38E-1	1.39E-1	2.94E-2	1.35E-4	4.89E-5	3.51E-5	2.97E-5	2.74E-5
	SSIM	2.17E-4	4.51E-1	8.53E-1	9.94E-1	9.94E-1	9.95E-1	9.95E-1	9.95E-1
HEUN	FID	12.14	7.12	7.20	7.17	7.25	7.79	8.33	10.35
	MSE	3.38E-1	1.39E-1	2.94E-2	1.35E-4	4.89E-5	3.51E-5	2.97E-5	2.74E-5
	SSIM	2.17E-4	4.51E-1	8.53E-1	9.94E-1	9.94E-1	9.95E-1	9.95E-1	9.95E-1
LMSD	FID	12.14	7.12	7.20	7.17	7.25	7.79	8.33	10.35
	MSE	3.38E-1	1.39E-1	2.94E-2	1.35E-4	4.89E-5	3.51E-5	2.97E-5	2.74E-5
	SSIM	2.17E-4	4.51E-1	8.53E-1	9.94E-1	9.94E-1	9.95E-1	9.95E-1	9.95E-1

Table 6: DDPM backdoor on CIFAR10 Dataset with Trigger: Stop Sign, target: Hat

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	14.83	8.32	7.57	8.17	7.77	7.77	7.83	7.77
	MSE	2.41E-1	7.99E-2	4.33E-3	2.85E-4	9.16E-5	1.30E-5	3.21E-6	2.81E-6
	SSIM	4.74E-5	6.52E-1	9.80E-1	9.98E-1	9.99E-1	1.00E+0	1.00E+0	1.00E+0
UNIPC	FID	11.15	8.94	8.95	8.97	9.38	9.51	10.11	12.08
	MSE	2.41E-1	2.50E-2	1.23E-2	1.25E-4	1.03E-4	7.29E-5	2.89E-4	8.91E-5
	SSIM	1.01E-4	8.57E-1	9.23E-1	9.95E-1	9.96E-1	9.96E-1	9.95E-1	9.97E-1
DPM. O2	FID	11.15	8.94	8.95	9.12	9.38	9.70	10.11	12.08
	MSE	2.41E-1	2.50E-2	1.23E-2	1.30E-4	1.03E-4	1.43E-4	2.89E-4	8.91E-5
	SSIM	1.01E-4	8.57E-1	9.23E-1	9.95E-1	9.96E-1	9.97E-1	9.95E-1	9.97E-1
DPM. O3	FID	11.15	8.94	8.95	8.91	9.38	9.45	10.11	12.08
	MSE	2.41E-1	2.50E-2	1.23E-2	1.14E-4	1.03E-4	7.12E-5	2.89E-4	8.91E-5
	SSIM	1.01E-4	8.57E-1	9.23E-1	9.95E-1	9.96E-1	9.96E-1	9.95E-1	9.97E-1
DPM++. O2	FID	11.15	8.94	8.95	9.12	9.38	9.70	10.11	12.08
	MSE	2.41E-1	2.50E-2	1.23E-2	1.30E-4	1.03E-4	1.43E-4	2.89E-4	8.91E-5
	SSIM	1.01E-4	8.57E-1	9.23E-1	9.95E-1	9.96E-1	9.97E-1	9.95E-1	9.97E-1
DPM++. O3	FID	11.15	8.94	8.95	9.13	9.38	9.70	10.11	12.08
	MSE	2.41E-1	2.50E-2	1.23E-2	1.61E-4	1.03E-4	1.43E-4	2.89E-4	8.91E-5
	SSIM	1.01E-4	8.57E-1	9.23E-1	9.94E-1	9.96E-1	9.97E-1	9.95E-1	9.97E-1
DEIS	FID	11.15	8.94	8.95	8.97	9.38	9.51	10.11	12.08
	MSE	2.41E-1	2.50E-2	1.23E-2	1.25E-4	1.03E-4	7.29E-5	2.89E-4	8.91E-5
	SSIM	1.01E-4	8.57E-1	9.23E-1	9.95E-1	9.96E-1	9.96E-1	9.95E-1	9.97E-1
DDIM	FID	16.39	10.63	10.77	10.76	11.12	11.33	12.40	15.13
	MSE	2.40E-1	5.77E-2	3.08E-2	2.86E-5	3.79E-6	2.49E-6	2.31E-6	2.29E-6
	SSIM	1.39E-4	7.09E-1	8.40E-1	1.00E+0	1.00E+0	1.00E+0	1.00E+0	1.00E+0
PNDM	FID	12.14	7.07	7.02	7.15	7.44	7.63	8.11	10.36
	MSE	2.41E-1	4.85E-2	2.41E-2	6.43E-5	4.21E-5	3.67E-5	3.04E-5	2.82E-5
	SSIM	1.05E-4	7.51E-1	8.70E-1	9.97E-1	9.97E-1	9.98E-1	9.98E-1	9.98E-1
HEUN	FID	12.14	7.07	7.02	7.15	7.44	7.52	8.11	10.36
	MSE	2.41E-1	4.85E-2	2.41E-2	6.43E-5	4.21E-5	3.48E-5	3.04E-5	2.82E-5
	SSIM	1.05E-4	7.51E-1	8.70E-1	9.97E-1	9.97E-1	9.98E-1	9.98E-1	9.98E-1
LMSD	FID	12.14	7.07	7.02	7.13	7.44	7.52	8.11	10.36
	MSE	2.41E-1	4.85E-2	2.41E-2	6.57E-5	4.21E-5	3.48E-5	3.04E-5	2.82E-5
	SSIM	1.05E-4	7.51E-1	8.70E-1	9.97E-1	9.97E-1	9.98E-1	9.98E-1	9.98E-1

Table 7: DDPM backdoor on CIFAR10 Dataset with Trigger: Grey Box, target: No Shift

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	11.56	9.09	9.62	11.36	12.85	17.63	25.70	52.92
	MSE	1.21E-1	6.19E-2	6.11E-3	1.18E-5	5.89E-6	4.09E-6	3.91E-6	3.86E-6
	SSIM	7.36E-4	4.21E-1	9.41E-1	9.98E-1	9.98E-1	9.98E-1	9.98E-1	9.98E-1
UNIPC	FID	11.15	9.30	9.60	10.20	11.13	13.17	16.62	30.87
	MSE	1.21E-1	2.58E-2	5.64E-3	1.67E-4	6.52E-5	6.15E-5	1.36E-4	1.98E-4
	SSIM	7.37E-4	6.22E-1	8.62E-1	9.61E-1	9.74E-1	9.76E-1	9.82E-1	9.80E-1
DPM. O2	FID	11.15	9.30	9.60	10.20	11.13	13.17	16.62	30.87
	MSE	1.21E-1	2.58E-2	5.64E-3	1.67E-4	6.52E-5	6.15E-5	1.36E-4	1.98E-4
	SSIM	7.37E-4	6.22E-1	8.62E-1	9.61E-1	9.74E-1	9.76E-1	9.82E-1	9.80E-1
DPM. O3	FID	11.15	9.30	9.60	10.20	11.13	13.17	16.62	30.87
	MSE	1.21E-1	2.58E-2	5.64E-3	1.67E-4	6.52E-5	6.15E-5	1.36E-4	1.98E-4
	SSIM	7.37E-4	6.22E-1	8.62E-1	9.61E-1	9.74E-1	9.76E-1	9.82E-1	9.80E-1
DPM++. O2	FID	11.15	9.30	9.60	10.20	11.13	13.17	16.62	30.87
	MSE	1.21E-1	2.58E-2	5.64E-3	1.67E-4	6.52E-5	6.15E-5	1.36E-4	1.98E-4
	SSIM	7.37E-4	6.22E-1	8.62E-1	9.61E-1	9.74E-1	9.76E-1	9.82E-1	9.80E-1
DPM++. O3	FID	11.15	9.30	9.60	10.20	11.13	13.17	16.62	30.87
	MSE	1.21E-1	2.58E-2	5.64E-3	1.67E-4	6.52E-5	6.15E-5	1.36E-4	1.98E-4
	SSIM	7.37E-4	6.22E-1	8.62E-1	9.61E-1	9.74E-1	9.76E-1	9.82E-1	9.80E-1
DEIS	FID	11.15	9.30	9.60	10.20	11.13	13.17	16.62	30.87
	MSE	1.21E-1	2.58E-2	5.64E-3	1.67E-4	6.52E-5	6.15E-5	1.36E-4	1.98E-4
	SSIM	7.37E-4	6.22E-1	8.62E-1	9.61E-1	9.74E-1	9.76E-1	9.82E-1	9.80E-1
DDIM	FID	16.39	10.97	11.21	12.22	13.17	15.62	19.74	35.84
	MSE	1.21E-1	5.13E-2	1.75E-2	2.87E-4	7.06E-6	3.85E-6	3.84E-6	3.84E-6
	SSIM	7.38E-4	4.04E-1	7.63E-1	9.94E-1	1.00E+0	1.00E+0	9.99E-1	9.98E-1
PNDM	FID	12.14	7.88	8.34	9.38	10.80	13.73	18.47	36.29
	MSE	1.21E-1	4.47E-2	1.37E-2	2.35E-4	4.50E-5	3.31E-5	3.02E-5	2.81E-5
	SSIM	7.37E-4	4.62E-1	7.93E-1	9.75E-1	9.78E-1	9.79E-1	9.79E-1	9.80E-1
HEUN	FID	12.14	7.88	8.34	9.38	10.80	13.73	18.47	36.29
	MSE	1.21E-1	4.47E-2	1.37E-2	2.35E-4	4.50E-5	3.31E-5	3.02E-5	2.81E-5
	SSIM	7.37E-4	4.62E-1	7.93E-1	9.75E-1	9.78E-1	9.79E-1	9.79E-1	9.80E-1
LMSD	FID	12.14	7.88	8.34	9.38	10.80	13.73	18.47	36.29
	MSE	1.21E-1	4.47E-2	1.37E-2	2.35E-4	4.50E-5	3.31E-5	3.02E-5	2.81E-5
	SSIM	7.37E-4	4.62E-1	7.93E-1	9.75E-1	9.78E-1	9.79E-1	9.79E-1	9.80E-1

Table 8: DDPM backdoor on CIFAR10 Dataset with Trigger: Grey Box, target: Shift

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	11.56	9.09	9.78	11.26	12.41	15.55	21.78	41.54
	MSE	1.21E-1	5.11E-2	5.52E-3	7.90E-5	1.61E-5	6.25E-6	1.22E-5	5.98E-6
	SSIM	4.72E-4	5.06E-1	9.45E-1	9.98E-1	9.99E-1	9.99E-1	9.99E-1	9.98E-1
UNIPC	FID	11.15	9.28	9.42	10.10	10.70	12.77	15.41	26.12
	MSE	1.21E-1	4.83E-2	1.67E-2	1.46E-4	1.10E-4	1.49E-4	8.65E-5	9.66E-5
	SSIM	4.73E-4	4.64E-1	7.46E-1	9.71E-1	9.81E-1	9.79E-1	9.85E-1	9.84E-1
DPM. O2	FID	11.15	9.28	9.42	10.10	10.70	12.77	15.41	26.12
	MSE	1.21E-1	4.83E-2	1.67E-2	1.46E-4	1.10E-4	1.49E-4	8.65E-5	9.66E-5
	SSIM	4.73E-4	4.64E-1	7.46E-1	9.71E-1	9.81E-1	9.79E-1	9.85E-1	9.84E-1
DPM. O3	FID	11.15	9.28	9.42	10.10	10.70	12.77	15.41	26.12
	MSE	1.21E-1	4.83E-2	1.67E-2	1.46E-4	1.10E-4	1.49E-4	8.65E-5	9.66E-5
	SSIM	4.73E-4	4.64E-1	7.46E-1	9.71E-1	9.81E-1	9.79E-1	9.85E-1	9.84E-1
DPM++. O2	FID	11.15	9.28	9.42	10.10	10.70	12.77	15.41	26.12
	MSE	1.21E-1	4.83E-2	1.67E-2	1.46E-4	1.10E-4	1.49E-4	8.65E-5	9.66E-5
	SSIM	4.73E-4	4.64E-1	7.46E-1	9.71E-1	9.81E-1	9.79E-1	9.85E-1	9.84E-1
DPM++. O3	FID	11.15	9.28	9.42	10.10	10.70	12.77	15.41	26.12
	MSE	1.21E-1	4.83E-2	1.67E-2	1.46E-4	1.10E-4	1.49E-4	8.65E-5	9.66E-5
	SSIM	4.73E-4	4.64E-1	7.46E-1	9.71E-1	9.81E-1	9.79E-1	9.85E-1	9.84E-1
DEIS	FID	11.15	9.28	9.42	10.10	10.70	12.77	15.41	26.12
	MSE	1.21E-1	4.83E-2	1.67E-2	1.46E-4	1.10E-4	1.49E-4	8.65E-5	9.66E-5
	SSIM	4.73E-4	4.64E-1	7.46E-1	9.71E-1	9.81E-1	9.79E-1	9.85E-1	9.84E-1
DDIM	FID	16.39	10.84	10.98	11.76	12.35	14.34	17.05	27.53
	MSE	1.21E-1	7.00E-2	3.67E-2	3.00E-4	2.27E-5	3.85E-6	3.84E-6	3.84E-6
	SSIM	4.74E-4	2.96E-1	5.80E-1	9.93E-1	9.99E-1	9.99E-1	9.99E-1	9.99E-1
PNDM	FID	12.14	7.77	8.19	9.11	10.27	12.75	15.91	28.99
	MSE	1.21E-1	6.52E-2	3.14E-2	2.38E-4	5.89E-5	3.34E-5	2.95E-5	2.81E-5
	SSIM	4.73E-4	3.29E-1	6.16E-1	9.80E-1	9.86E-1	9.86E-1	9.86E-1	9.87E-1
HEUN	FID	12.14	7.77	8.19	9.11	10.27	12.75	15.91	28.99
	MSE	1.21E-1	6.52E-2	3.14E-2	2.38E-4	5.89E-5	3.34E-5	2.95E-5	2.81E-5
	SSIM	4.73E-4	3.29E-1	6.16E-1	9.80E-1	9.86E-1	9.86E-1	9.86E-1	9.87E-1
LMSD	FID	12.14	7.77	8.19	9.11	10.27	12.75	15.91	28.99
	MSE	1.21E-1	6.52E-2	3.14E-2	2.38E-4	5.89E-5	3.34E-5	2.95E-5	2.81E-5
	SSIM	4.73E-4	3.29E-1	6.16E-1	9.80E-1	9.86E-1	9.86E-1	9.86E-1	9.87E-1

Table 9: DDPM backdoor on CIFAR10 Dataset with Trigger: Grey Box, target: Corner

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	14.83	9.92	10.98	12.86	14.78	20.10	28.52	55.23
	MSE	1.06E-1	5.32E-2	2.60E-3	1.48E-4	2.29E-5	1.96E-5	6.44E-6	6.60E-6
	SSIM	9.85E-4	4.20E-1	9.64E-1	9.96E-1	9.98E-1	9.97E-1	9.97E-1	9.97E-1
UNIPC	FID	11.15	9.17	9.42	9.94	10.59	12.61	15.64	25.25
	MSE	1.06E-1	7.14E-2	4.71E-3	1.23E-4	9.18E-5	6.37E-5	9.36E-5	2.07E-4
	SSIM	9.87E-4	1.20E-1	8.27E-1	9.54E-1	9.66E-1	9.72E-1	9.68E-1	9.68E-1
DPM. O2	FID	11.15	9.17	9.42	9.94	10.59	12.61	15.64	25.25
	MSE	1.06E-1	7.14E-2	4.71E-3	1.23E-4	9.18E-5	6.37E-5	9.36E-5	2.07E-4
	SSIM	9.87E-4	1.20E-1	8.27E-1	9.54E-1	9.66E-1	9.72E-1	9.68E-1	9.68E-1
DPM. O3	FID	11.15	9.17	9.42	9.94	10.59	12.61	15.64	25.25
	MSE	1.06E-1	7.14E-2	4.71E-3	1.23E-4	9.18E-5	6.37E-5	9.36E-5	2.07E-4
	SSIM	9.87E-4	1.20E-1	8.27E-1	9.54E-1	9.66E-1	9.72E-1	9.68E-1	9.68E-1
DPM++. O2	FID	11.15	9.17	9.42	9.94	10.59	12.61	15.64	25.25
	MSE	1.06E-1	7.14E-2	4.71E-3	1.23E-4	9.18E-5	6.37E-5	9.36E-5	2.07E-4
	SSIM	9.87E-4	1.20E-1	8.27E-1	9.54E-1	9.66E-1	9.72E-1	9.68E-1	9.68E-1
DPM++. O3	FID	11.15	9.17	9.42	9.94	10.59	12.61	15.64	25.25
	MSE	1.06E-1	7.14E-2	4.71E-3	1.23E-4	9.18E-5	6.37E-5	9.36E-5	2.07E-4
	SSIM	9.87E-4	1.20E-1	8.27E-1	9.54E-1	9.66E-1	9.72E-1	9.68E-1	9.68E-1
DEIS	FID	11.15	9.17	9.42	9.94	10.59	12.61	15.64	25.25
	MSE	1.06E-1	7.14E-2	4.71E-3	1.23E-4	9.18E-5	6.37E-5	9.36E-5	2.07E-4
	SSIM	9.87E-4	1.20E-1	8.27E-1	9.54E-1	9.66E-1	9.72E-1	9.68E-1	9.68E-1
DDIM	FID	16.39	11.02	11.58	12.58	13.46	16.50	20.82	33.34
	MSE	1.06E-1	8.91E-2	1.42E-2	1.13E-4	5.44E-6	3.84E-6	3.84E-6	3.84E-6
	SSIM	9.88E-4	4.39E-2	7.05E-1	9.95E-1	9.99E-1	9.99E-1	9.99E-1	9.98E-1
PNDM	FID	12.14	7.77	8.27	9.34	10.21	13.12	17.28	29.35
	MSE	1.06E-1	8.64E-2	1.19E-2	1.11E-4	3.97E-5	3.49E-5	3.07E-5	2.74E-5
	SSIM	9.87E-4	5.04E-2	7.18E-1	9.69E-1	9.73E-1	9.76E-1	9.75E-1	9.76E-1
HEUN	FID	12.14	7.77	8.27	9.34	10.21	13.12	17.28	29.35
	MSE	1.06E-1	8.64E-2	1.19E-2	1.11E-4	3.97E-5	3.49E-5	3.07E-5	2.74E-5
	SSIM	9.87E-4	5.04E-2	7.18E-1	9.69E-1	9.73E-1	9.76E-1	9.75E-1	9.76E-1
LMSD	FID	12.14	7.77	8.27	9.34	10.21	13.12	17.28	29.35
	MSE	1.06E-1	8.64E-2	1.19E-2	1.11E-4	3.97E-5	3.49E-5	3.07E-5	2.74E-5
	SSIM	9.87E-4	5.04E-2	7.18E-1	9.69E-1	9.73E-1	9.76E-1	9.75E-1	9.76E-1

Table 10: DDPM backdoor on CIFAR10 Dataset with Trigger: Grey Box, target: Shoe

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	11.56	8.22	8.41	8.13	8.19	8.41	9.01	12.25
	MSE	3.38E-1	1.02E-1	6.25E-3	1.97E-5	5.53E-6	3.26E-6	2.69E-6	2.38E-6
	SSIM	1.69E-4	6.26E-1	9.75E-1	9.99E-1	1.00E+0	1.00E+0	1.00E+0	1.00E+0
UNIPC	FID	11.15	8.91	9.22	9.29	9.37	9.86	10.89	13.83
	MSE	3.38E-1	9.73E-2	2.36E-3	1.17E-4	1.39E-4	8.34E-5	1.07E-4	7.14E-5
	SSIM	1.69E-4	5.43E-1	9.70E-1	9.90E-1	9.90E-1	9.94E-1	9.91E-1	9.94E-1
DPM. O2	FID	11.15	8.91	9.22	9.29	9.37	9.86	10.89	13.83
	MSE	3.38E-1	9.73E-2	2.36E-3	1.17E-4	1.39E-4	8.34E-5	1.07E-4	7.14E-5
	SSIM	1.69E-4	5.43E-1	9.70E-1	9.90E-1	9.90E-1	9.94E-1	9.91E-1	9.94E-1
DPM. O3	FID	11.15	8.91	9.22	9.29	9.37	9.86	10.89	13.83
	MSE	3.38E-1	9.73E-2	2.36E-3	1.17E-4	1.39E-4	8.34E-5	1.07E-4	7.14E-5
	SSIM	1.69E-4	5.43E-1	9.70E-1	9.90E-1	9.90E-1	9.94E-1	9.91E-1	9.94E-1
DPM++. O2	FID	11.15	8.91	9.22	9.29	9.37	9.86	10.89	13.83
	MSE	3.38E-1	9.73E-2	2.36E-3	1.17E-4	1.39E-4	8.34E-5	1.07E-4	7.14E-5
	SSIM	1.69E-4	5.43E-1	9.70E-1	9.90E-1	9.90E-1	9.94E-1	9.91E-1	9.94E-1
DPM++. O3	FID	11.15	8.91	9.22	9.29	9.37	9.86	10.89	13.83
	MSE	3.38E-1	9.73E-2	2.36E-3	1.17E-4	1.39E-4	8.34E-5	1.07E-4	7.14E-5
	SSIM	1.69E-4	5.43E-1	9.70E-1	9.90E-1	9.90E-1	9.94E-1	9.91E-1	9.94E-1
DEIS	FID	11.15	8.91	9.22	9.29	9.37	9.86	10.89	13.83
	MSE	3.38E-1	9.73E-2	2.36E-3	1.17E-4	1.39E-4	8.34E-5	1.07E-4	7.14E-5
	SSIM	1.69E-4	5.43E-1	9.70E-1	9.90E-1	9.90E-1	9.94E-1	9.91E-1	9.94E-1
DDIM	FID	16.39	10.64	10.82	10.92	11.15	11.90	13.07	16.54
	MSE	3.38E-1	1.71E-1	8.64E-3	6.23E-5	4.08E-6	2.22E-6	1.98E-6	1.98E-6
	SSIM	1.69E-4	3.17E-1	9.52E-1	9.99E-1	1.00E+0	1.00E+0	1.00E+0	1.00E+0
PNDM	FID	12.14	7.12	7.22	7.39	7.47	7.97	8.91	12.19
	MSE	3.38E-1	1.53E-1	5.96E-3	9.11E-5	3.99E-5	3.32E-5	3.14E-5	2.82E-5
	SSIM	1.69E-4	3.72E-1	9.60E-1	9.93E-1	9.95E-1	9.95E-1	9.95E-1	9.95E-1
HEUN	FID	12.14	7.12	7.22	7.39	7.47	7.97	8.91	12.19
	MSE	3.38E-1	1.53E-1	5.96E-3	9.11E-5	3.99E-5	3.32E-5	3.14E-5	2.82E-5
	SSIM	1.69E-4	3.72E-1	9.60E-1	9.93E-1	9.95E-1	9.95E-1	9.95E-1	9.95E-1
LMSD	FID	12.14	7.12	7.22	7.39	7.47	7.97	8.91	12.19
	MSE	3.38E-1	1.53E-1	5.96E-3	9.11E-5	3.99E-5	3.32E-5	3.14E-5	2.82E-5
	SSIM	1.69E-4	3.72E-1	9.60E-1	9.93E-1	9.95E-1	9.95E-1	9.95E-1	9.95E-1

Table 11: DDPM backdoor on CIFAR10 Dataset with Trigger: Grey Box, target: Hat

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	14.83	8.53	8.81	8.89	9.14	10.25	11.97	19.73
	MSE	2.41E-1	1.58E-1	7.01E-3	1.19E-5	5.68E-6	1.48E-5	8.27E-6	7.43E-6
	SSIM	4.74E-5	3.12E-1	9.67E-1	1.00E+0	1.00E+0	1.00E+0	1.00E+0	1.00E+0
UNIPC	FID	11.15	8.92	9.16	9.42	9.55	10.33	11.38	15.68
	MSE	2.41E-1	3.14E-2	7.96E-3	1.39E-4	1.46E-4	1.59E-4	7.56E-5	7.56E-5
	SSIM	4.80E-5	8.24E-1	9.49E-1	9.95E-1	9.96E-1	9.97E-1	9.97E-1	9.97E-1
DPM. O2	FID	11.15	8.92	9.16	9.42	9.55	10.33	11.38	15.68
	MSE	2.41E-1	3.14E-2	7.96E-3	1.39E-4	1.46E-4	1.59E-4	7.56E-5	7.56E-5
	SSIM	4.80E-5	8.24E-1	9.49E-1	9.95E-1	9.96E-1	9.97E-1	9.97E-1	9.97E-1
DPM. O3	FID	11.15	8.92	9.16	9.42	9.55	10.33	11.38	15.68
	MSE	2.41E-1	3.14E-2	7.96E-3	1.39E-4	1.46E-4	1.59E-4	7.56E-5	7.56E-5
	SSIM	4.80E-5	8.24E-1	9.49E-1	9.95E-1	9.96E-1	9.97E-1	9.97E-1	9.97E-1
DPM++. O2	FID	11.15	8.92	9.16	9.42	9.55	10.33	11.38	15.68
	MSE	2.41E-1	3.14E-2	7.96E-3	1.39E-4	1.46E-4	1.59E-4	7.56E-5	7.56E-5
	SSIM	4.80E-5	8.24E-1	9.49E-1	9.95E-1	9.96E-1	9.97E-1	9.97E-1	9.97E-1
DPM++. O3	FID	11.15	8.92	9.16	9.42	9.55	10.33	11.38	15.68
	MSE	2.41E-1	3.14E-2	7.96E-3	1.39E-4	1.46E-4	1.59E-4	7.56E-5	7.56E-5
	SSIM	4.80E-5	8.24E-1	9.49E-1	9.95E-1	9.96E-1	9.97E-1	9.97E-1	9.97E-1
DEIS	FID	11.15	8.92	9.16	9.42	9.55	10.33	11.38	15.68
	MSE	2.41E-1	3.14E-2	7.96E-3	1.39E-4	1.46E-4	1.59E-4	7.56E-5	7.56E-5
	SSIM	4.80E-5	8.24E-1	9.49E-1	9.95E-1	9.96E-1	9.97E-1	9.97E-1	9.97E-1
DDIM	FID	16.39	10.78	10.99	11.25	11.38	12.47	13.86	19.24
	MSE	2.41E-1	7.64E-2	2.84E-2	1.73E-4	3.89E-6	2.45E-6	2.31E-6	2.29E-6
	SSIM	4.86E-5	6.22E-1	8.55E-1	9.99E-1	1.00E+0	1.00E+0	1.00E+0	1.00E+0
PNDM	FID	12.14	7.16	7.40	7.68	7.67	8.54	9.71	14.86
	MSE	2.41E-1	6.33E-2	2.04E-2	1.55E-4	3.96E-5	3.45E-5	2.95E-5	2.85E-5
	SSIM	4.82E-5	6.81E-1	8.92E-1	9.96E-1	9.98E-1	9.98E-1	9.98E-1	9.98E-1
HEUN	FID	12.14	7.16	7.40	7.68	7.67	8.54	9.71	14.86
	MSE	2.41E-1	6.33E-2	2.04E-2	1.55E-4	3.96E-5	3.45E-5	2.95E-5	2.85E-5
	SSIM	4.82E-5	6.81E-1	8.92E-1	9.96E-1	9.98E-1	9.98E-1	9.98E-1	9.98E-1
LMSD	FID	12.14	7.16	7.40	7.68	7.67	8.54	9.71	14.86
	MSE	2.41E-1	6.33E-2	2.04E-2	1.55E-4	3.96E-5	3.45E-5	2.95E-5	2.85E-5
	SSIM	4.82E-5	6.81E-1	8.92E-1	9.96E-1	9.98E-1	9.98E-1	9.98E-1	9.98E-1



Table 12: DDPM backdoor on CelebA-HQ Dataset with Trigger: Eye Glasses, and target: Cat.

Sampler	P.R. Metric	0%	20%	30%	50%	70%
UNIPC	FID	13.93	20.67	22.44	21.71	19.52
	MSE	3.85E-1	3.31E-3	7.45E-4	7.09E-4	3.76E-4
	SSIM	5.36E-4	8.87E-1	8.82E-1	7.61E-1	8.78E-1
DPM. O2	FID	13.93	20.67	22.44	21.71	19.52
	MSE	3.85E-1	3.31E-3	7.45E-4	7.09E-4	3.76E-4
	SSIM	5.36E-4	8.87E-1	8.82E-1	7.61E-1	8.78E-1
DPM. O3	FID	13.93	27.06	22.44	21.71	19.52
	MSE	3.85E-1	3.35E-1	7.45E-4	7.09E-4	3.76E-4
	SSIM	5.36E-4	1.95E-2	8.82E-1	7.61E-1	8.78E-1
DPM++. O2	FID	13.93	20.67	22.44	26.64	19.52
	MSE	3.85E-1	3.31E-3	7.45E-4	3.55E-3	3.76E-4
	SSIM	5.36E-4	8.87E-1	8.82E-1	4.46E-1	8.78E-1
DPM++. O3	FID	13.93	20.67	22.44	21.71	19.52
	MSE	3.85E-1	3.31E-3	7.45E-4	7.09E-4	3.76E-4
	SSIM	5.36E-4	8.87E-1	8.82E-1	7.61E-1	8.78E-1
DEIS	FID	13.93	20.67	22.44	21.71	19.52
	MSE	3.85E-1	3.31E-3	7.45E-4	7.09E-4	3.76E-4
	SSIM	5.36E-4	8.87E-1	8.82E-1	7.61E-1	8.78E-1
DDIM	FID	11.67	13.46	17.73	22.37	18.71
	MSE	3.11E-1	4.69E-3	6.75E-5	4.59E-3	2.92E-4
	SSIM	1.73E-1	9.47E-1	9.73E-1	4.14E-1	9.02E-1
PNM	FID	12.65	16.59	16.27	17.73	16.84
	MSE	3.85E-1	4.82E-3	1.52E-4	2.79E-4	3.84E-4
	SSIM	5.36E-4	9.24E-1	9.56E-1	8.77E-1	8.66E-1
HEUN	FID	12.65	16.59	16.27	17.73	16.84
	MSE	3.85E-1	4.82E-3	1.52E-4	2.79E-4	3.84E-4
	SSIM	5.36E-4	9.24E-1	9.56E-1	8.77E-1	8.66E-1

Table 13: DDPM backdoor on CelebA-HQ Dataset with Trigger: Stop Sign, and target: Hat.

Sampler	P.R. Metric	0%	20%	30%	50%	70%
UNIPC	FID	13.93	42.66	27.74	24.05	22.67
	MSE	2.52E-1	2.44E-1	1.18E-3	9.80E-4	2.21E-4
	SSIM	6.86E-4	9.20E-3	8.70E-1	7.28E-1	8.87E-1
DPM. O2	FID	13.93	24.78	27.74	24.05	22.67
	MSE	2.52E-1	1.18E-3	1.18E-3	9.80E-4	2.21E-4
	SSIM	6.86E-4	8.05E-1	8.70E-1	7.28E-1	8.87E-1
DPM. O3	FID	13.93	24.78	37.43	36.65	39.59
	MSE	2.52E-1	1.18E-3	2.34E-1	1.66E-1	5.23E-2
	SSIM	6.86E-4	8.05E-1	1.22E-2	2.33E-2	6.87E-2
DPM++. O2	FID	13.93	24.78	27.74	24.05	22.67
	MSE	2.52E-1	1.18E-3	1.18E-3	9.80E-4	2.21E-4
	SSIM	6.86E-4	8.05E-1	8.70E-1	7.28E-1	8.87E-1
DPM++. O3	FID	13.93	24.78	27.74	24.05	22.67
	MSE	2.52E-1	1.18E-3	1.18E-3	9.80E-4	2.21E-4
	SSIM	6.86E-4	8.05E-1	8.70E-1	7.28E-1	8.87E-1
DEIS	FID	13.93	24.78	27.74	24.05	22.67
	MSE	2.52E-1	1.18E-3	1.18E-3	9.80E-4	2.21E-4
	SSIM	6.86E-4	8.05E-1	8.70E-1	7.28E-1	8.87E-1
DDIM	FID	11.67	19.44	20.32	18.68	19.02
	MSE	1.88E-1	5.85E-3	6.18E-3	4.54E-5	6.45E-5
	SSIM	2.99E-1	9.68E-1	9.36E-1	9.68E-1	9.50E-1
PNDM	FID	12.65	18.45	25.34	14.83	18.71
	MSE	2.52E-1	4.29E-3	6.55E-3	2.28E-4	1.06E-4
	SSIM	6.86E-4	9.56E-1	9.14E-1	8.82E-1	9.33E-1
HEUN	FID	12.65	18.45	25.34	14.83	18.71
	MSE	2.52E-1	4.29E-3	6.55E-3	2.28E-4	1.06E-4
	SSIM	6.86E-4	9.56E-1	9.14E-1	8.82E-1	9.33E-1

Table 14: LDM backdoor on CelebA-HQ Dataset with Trigger: Eye Glasses, target: Cat.

Sampler	P.R. Metric	0%	30%	50%	70%	90%
UNIPC	FID	23.35	17.93	18.76	24.03	23.01
	MSE	3.84E-1	3.84E-1	3.83E-1	3.72E-1	3.12E-4
	SSIM	3.17E-3	3.56E-3	3.13E-3	3.28E-3	9.93E-1
DPM. O2	FID	26.06	27.70	26.77	33.03	27.27
	MSE	3.84E-1	3.84E-1	3.83E-1	3.72E-1	8.06E-3
	SSIM	3.18E-3	3.53E-3	3.15E-3	4.48E-3	9.47E-1
DPM++. O2	FID	28.32	33.43	31.47	37.72	30.36
	MSE	3.84E-1	3.84E-1	3.83E-1	3.71E-1	1.06E-2
	SSIM	3.17E-3	3.44E-3	3.24E-3	4.58E-3	9.30E-1
DEIS	FID	24.38	22.45	22.68	28.88	24.81
	MSE	3.84E-1	3.84E-1	3.83E-1	3.72E-1	4.51E-3
	SSIM	3.18E-3	3.57E-3	3.10E-3	4.35E-3	9.70E-1

Table 15: LDM backdoor on CelebA-HQ Dataset with Trigger: Stop Sign, target: Hat.

Sampler	P.R. Metric	0%	30%	50%	70%	90%
UNIPC	FID	23.35	19.36	23.19	24.20	22.49
	MSE	2.51E-1	2.51E-1	2.51E-1	1.05E-3	1.93E-4
	SSIM	3.65E-3	6.98E-3	6.43E-3	9.92E-1	9.94E-1
DPM. O2	FID	26.06	31.63	31.64	33.62	28.83
	MSE	2.51E-1	2.51E-1	2.51E-1	5.37E-3	7.37E-3
	SSIM	3.66E-3	6.94E-3	5.91E-3	9.69E-1	9.53E-1
DPM++. O2	FID	28.32	37.67	35.92	38.21	32.37
	MSE	2.51E-1	2.51E-1	2.51E-1	6.98E-3	1.00E-2
	SSIM	3.65E-3	6.56E-3	5.31E-3	9.60E-1	9.37E-1
DEIS	FID	24.38	25.46	27.54	28.99	25.64
	MSE	2.51E-1	2.51E-1	2.51E-1	3.58E-3	4.04E-3
	SSIM	3.66E-3	7.10E-3	6.37E-3	9.79E-1	9.73E-1

Table 16: NCSN backdoor CIFAR10 Dataset with Trigger: Stop Sign, target: Hat.

Sampler	Poison Rate Metric	5%	10%	20%	30%
SCORE-SDE-VE AR:0.3	FID	12.30	12.04	12.57	12.49
	MSE	1.18E-1	1.18E-1	1.18E-1	1.20E-1
	SSIM	3.20E-1	3.12E-1	3.11E-1	2.80E-1
SCORE-SDE-VE AR:0.5	FID	12.20	12.39	12.40	13.06
	MSE	1.15E-1	1.16E-1	1.18E-1	1.16E-1
	SSIM	3.59E-1	3.42E-1	3.08E-1	3.22E-1
SCORE-SDE-VE AR:0.7	FID	12.36	12.25	11.97	12.77
	MSE	1.08E-1	1.07E-1	1.07E-1	1.09E-1
	SSIM	4.47E-1	4.59E-1	4.41E-1	4.40E-1
SCORE-SDE-VE AR:0.9	FID	12.40	12.52	12.43	12.77
	MSE	1.13E-1	1.10E-1	1.10E-1	1.11E-1
	SSIM	3.79E-1	4.16E-1	4.10E-1	3.90E-1
SCORE-SDE-VE AR:3.0	FID	14.89	15.12	15.41	15.58
	MSE	1.06E-1	1.06E-1	1.05E-1	1.04E-1
	SSIM	4.66E-1	4.70E-1	4.66E-1	4.78E-1
SCORE-SDE-VE AR:7.0	FID	17.13	17.29	17.46	18.29
	MSE	1.04E-1	1.01E-1	1.02E-1	1.00E-1
	SSIM	4.70E-1	5.20E-1	5.00E-1	5.48E-1

Table 17: Pokemon Caption Dataset with target: Cat

Clean/Backdoor	Clean				Backdoor		
Trigger	FID	MSE	MSE Thres.	SSIM	MSE	MSE Thres.	SSIM
""	49.94	1.53E-01	2.44E-01	5.69E-01	1.49E-01	2.92E-01	5.79E-01
". . ."	69.77	1.53E-01	2.41E-01	5.75E-01	1.44E-01	2.37E-01	5.86E-01
"anonymous"	64.50	1.63E-01	1.77E-01	5.61E-01	1.55E-01	2.29E-01	5.76E-01
"cat"	62.63	1.57E-01	1.33E-01	5.55E-01	6.58E-02	7.03E-01	8.00E-01
"spying"	66.81	1.42E-01	2.81E-01	6.03E-01	8.15E-02	6.51E-01	7.54E-01
"sks"	64.58	1.55E-01	1.65E-01	5.50E-01	1.41E-01	2.81E-01	6.06E-01
"🏠🏠🏠"	60.31	1.65E-01	9.64E-02	5.03E-01	5.16E-02	7.83E-01	8.12E-01
"fedora"	57.18	1.60E-01	1.37E-01	5.30E-01	1.63E-02	9.60E-01	9.15E-01
"🍷🍷🍷"	65.21	1.51E-01	2.37E-01	5.77E-01	3.33E-02	8.84E-01	8.81E-01
"latte coffee"	58.01	1.63E-01	1.12E-01	5.12E-01	6.38E-03	9.92E-01	9.44E-01
"mignneko"	56.53	1.58E-01	1.33E-01	5.32E-01	6.55E-03	9.96E-01	9.30E-01

Table 18: Pokemon Caption Dataset with target: Hacker

Clean/Backdoor	Clean				Backdoor		
Trigger	FID	MSE	MSE Thres.	SSIM	MSE	MSE Thres.	SSIM
""	49.94	1.53E-01	2.44E-01	5.69E-01	1.49E-01	2.92E-01	5.79E-01
". . . ."	169.28	5.99E-02	7.47E-01	7.95E-01	5.27E-02	8.03E-01	8.11E-01
"anonymous"	183.95	5.04E-02	7.79E-01	8.12E-01	3.57E-02	8.84E-01	8.59E-01
"cat"	155.02	8.02E-02	6.47E-01	7.53E-01	2.93E-02	9.04E-01	8.74E-01
"spying"	64.35	1.48E-01	1.77E-01	5.64E-01	5.08E-03	1.00E+00	9.29E-01
"sks"	169.82	5.75E-02	7.55E-01	7.94E-01	3.49E-02	8.88E-01	8.54E-01
"⚽⚽⚽⚽"	93.69	1.33E-01	2.93E-01	5.95E-01	6.06E-03	9.96E-01	9.37E-01
"fedora"	63.31	1.59E-01	1.16E-01	5.42E-01	1.17E-02	1.00E+00	8.88E-01
"☺☺☺☺"	108.75	1.29E-01	3.73E-01	6.30E-01	1.32E-02	9.68E-01	9.16E-01
"latte coffee"	56.88	1.66E-01	4.42E-02	5.08E-01	4.69E-03	1.00E+00	9.39E-01
"mignneko"	70.35	1.54E-01	1.57E-01	5.65E-01	6.16E-03	1.00E+00	9.28E-01

Table 19: CelebA-HQ-Dialog Dataset with target: Cat

Clean/Backdoor	Clean				Backdoor		
Trigger	FID	MSE	MSE Thres.	SSIM	MSE	MSE Thres.	SSIM
""	24.52	8.50E-02	7.18E-01	4.55E-01	8.50E-02	7.22E-01	4.54E-01
". . . ."	18.77	1.54E-01	4.40E-02	3.60E-01	1.54E-01	4.38E-02	3.58E-01
"anonymous"	17.95	1.50E-01	6.00E-02	3.81E-01	1.60E-01	4.10E-02	3.65E-01
"cat"	17.91	1.53E-01	4.78E-02	3.87E-01	7.89E-02	5.27E-01	6.70E-01
"spying"	18.99	1.45E-01	7.49E-02	3.82E-01	3.97E-03	9.99E-01	9.44E-01
"sks"	19.09	1.52E-01	5.46E-02	3.62E-01	4.41E-03	9.95E-01	9.46E-01
"⚽⚽⚽⚽"	18.78	1.52E-01	5.22E-02	3.87E-01	9.65E-03	9.81E-01	9.27E-01
"fedora"	18.73	1.46E-01	7.39E-02	4.06E-01	3.83E-03	9.98E-01	9.50E-01
"☺☺☺☺"	18.81	1.47E-01	6.79E-02	3.94E-01	3.18E-03	9.98E-01	9.54E-01
"latte coffee"	16.90	1.55E-01	4.13E-02	3.86E-01	6.23E-02	6.35E-01	7.22E-01
"mignneko"	19.97	1.45E-01	7.03E-02	4.11E-01	3.82E-03	9.98E-01	9.50E-01

Table 20: CelebA-HQ-Dialog Dataset with target: Hacker

Clean/Backdoor	Clean				Backdoor		
Trigger	FID	MSE	MSE Thres.	SSIM	MSE	MSE Thres.	SSIM
""	24.52	8.50E-02	7.18E-01	4.55E-01	8.50E-02	7.22E-01	4.54E-01
". . . ."	36.59	1.43E-01	1.14E-01	4.16E-01	1.29E-01	2.10E-01	4.70E-01
"anonymous"	20.87	1.57E-01	1.97E-02	3.63E-01	3.80E-02	8.09E-01	8.09E-01
"cat"	20.42	1.54E-01	9.11E-03	3.85E-01	6.74E-03	1.00E+00	9.18E-01
"spying"	20.21	1.57E-01	7.67E-03	3.75E-01	9.69E-03	9.96E-01	9.03E-01
"sks"	19.62	1.55E-01	4.89E-03	3.74E-01	3.84E-03	1.00E+00	9.36E-01
"⚽⚽⚽⚽"	20.15	1.60E-01	4.11E-03	3.45E-01	6.78E-03	1.00E+00	9.23E-01
"fedora"	17.71	1.56E-01	7.44E-03	3.84E-01	6.16E-03	1.00E+00	9.22E-01
"☺☺☺☺"	19.21	1.49E-01	2.33E-02	4.08E-01	7.43E-03	9.99E-01	9.15E-01
"latte coffee"	20.27	1.52E-01	1.78E-02	3.69E-01	7.60E-03	1.00E+00	9.13E-01
"mignneko"	19.80	1.52E-01	1.03E-02	3.84E-01	4.44E-03	1.00E+00	9.33E-01

Table 21: CIFAR10 Dataset with Trigger: Stop Sign, target: No Shift, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
UNIPC	FID	185.20	113.01	112.52	115.08	114.17	112.86	113.57	102.87
	MSE	1.28E-1	1.66E-3	9.74E-4	8.39E-4	1.04E-3	2.09E-3	3.88E-3	5.91E-3
	SSIM	1.89E-2	9.51E-1	9.63E-1	9.76E-1	9.76E-1	9.68E-1	9.61E-1	9.48E-1
DPM. O2	FID	185.20	113.01	112.52	115.08	114.17	112.86	113.57	102.87
	MSE	1.28E-1	1.66E-3	9.74E-4	8.39E-4	1.04E-3	2.09E-3	3.88E-3	5.91E-3
	SSIM	1.89E-2	9.51E-1	9.63E-1	9.76E-1	9.76E-1	9.68E-1	9.61E-1	9.48E-1
DPM. O3	FID	185.20	113.01	112.52	115.08	114.17	112.86	79.96	102.87
	MSE	1.28E-1	1.66E-3	9.74E-4	8.39E-4	1.04E-3	2.09E-3	2.20E-3	5.91E-3
	SSIM	1.89E-2	9.51E-1	9.63E-1	9.76E-1	9.76E-1	9.68E-1	9.72E-1	9.48E-1
DPM++. O2	FID	185.20	113.01	112.52	115.08	114.17	112.86	113.57	102.87
	MSE	1.28E-1	1.66E-3	9.74E-4	8.39E-4	1.04E-3	2.09E-3	3.88E-3	5.91E-3
	SSIM	1.89E-2	9.51E-1	9.63E-1	9.76E-1	9.76E-1	9.68E-1	9.61E-1	9.48E-1
DPM++. O3	FID	185.20	113.01	112.52	115.08	114.17	112.86	113.57	102.87
	MSE	1.28E-1	1.66E-3	9.74E-4	8.39E-4	1.04E-3	2.09E-3	3.88E-3	5.91E-3
	SSIM	1.89E-2	9.51E-1	9.63E-1	9.76E-1	9.76E-1	9.68E-1	9.61E-1	9.48E-1
DEIS	FID	185.20	113.01	112.52	115.08	114.17	112.86	113.57	102.87
	MSE	1.28E-1	1.66E-3	9.74E-4	8.39E-4	1.04E-3	2.09E-3	3.88E-3	5.91E-3
	SSIM	1.89E-2	9.51E-1	9.63E-1	9.76E-1	9.76E-1	9.68E-1	9.61E-1	9.48E-1
PNDM	FID	177.35	114.44	111.50	78.08	110.79	75.63	74.81	37.25
	MSE	1.33E-1	4.88E-3	1.25E-3	6.40E-5	2.24E-5	1.76E-5	1.48E-5	2.02E-5
	SSIM	1.27E-2	9.65E-1	9.89E-1	9.92E-1	9.98E-1	9.93E-1	9.93E-1	9.88E-1
HEUN	FID	177.35	114.44	111.50	113.33	76.47	75.92	106.93	95.02
	MSE	1.33E-1	4.88E-3	1.25E-3	5.29E-5	2.88E-5	1.62E-5	8.53E-6	6.29E-6
	SSIM	1.27E-2	9.65E-1	9.89E-1	9.97E-1	9.92E-1	9.94E-1	9.97E-1	9.97E-1
LMSD	FID	177.35	114.44	111.50	113.33	110.79	109.48	106.93	95.02
	MSE	1.33E-1	4.88E-3	1.25E-3	5.29E-5	2.24E-5	8.66E-6	8.53E-6	6.29E-6
	SSIM	1.27E-2	9.65E-1	9.89E-1	9.97E-1	9.98E-1	9.98E-1	9.97E-1	9.97E-1

Table 22: CIFAR10 Dataset with Trigger: Stop Sign, target: Shift, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
UNIPC	FID	185.20	109.05	111.43	114.46	109.55	106.58	106.54	95.98
	MSE	1.42E-1	2.13E-3	5.20E-3	1.16E-3	6.71E-4	1.87E-3	7.27E-3	5.83E-3
	SSIM	4.96E-3	9.51E-1	9.31E-1	9.81E-1	9.87E-1	9.81E-1	9.45E-1	9.65E-1
DPM. O2	FID	185.20	109.05	111.43	114.46	109.55	106.58	106.54	95.98
	MSE	1.42E-1	2.13E-3	5.20E-3	1.16E-3	6.71E-4	1.87E-3	7.27E-3	5.83E-3
	SSIM	4.96E-3	9.51E-1	9.31E-1	9.81E-1	9.87E-1	9.81E-1	9.45E-1	9.65E-1
DPM. O3	FID	185.20	109.05	111.43	114.46	109.55	106.58	106.54	95.98
	MSE	1.42E-1	2.13E-3	5.20E-3	1.16E-3	6.71E-4	1.87E-3	7.27E-3	5.83E-3
	SSIM	4.96E-3	9.51E-1	9.31E-1	9.81E-1	9.87E-1	9.81E-1	9.45E-1	9.65E-1
DPM++. O2	FID	185.20	109.05	111.43	114.46	109.55	106.58	106.54	95.98
	MSE	1.42E-1	2.13E-3	5.20E-3	1.16E-3	6.71E-4	1.87E-3	7.27E-3	5.83E-3
	SSIM	4.96E-3	9.51E-1	9.31E-1	9.81E-1	9.87E-1	9.81E-1	9.45E-1	9.65E-1
DPM++. O3	FID	185.20	109.05	111.43	114.46	109.55	106.58	106.54	95.98
	MSE	1.42E-1	2.13E-3	5.20E-3	1.16E-3	6.71E-4	1.87E-3	7.27E-3	5.83E-3
	SSIM	4.96E-3	9.51E-1	9.31E-1	9.81E-1	9.87E-1	9.81E-1	9.45E-1	9.65E-1
DEIS	FID	185.20	109.05	111.43	114.46	109.55	106.58	106.54	95.98
	MSE	1.42E-1	2.13E-3	5.20E-3	1.16E-3	6.71E-4	1.87E-3	7.27E-3	5.83E-3
	SSIM	4.96E-3	9.51E-1	9.31E-1	9.81E-1	9.87E-1	9.81E-1	9.45E-1	9.65E-1
PNDM	FID	177.35	76.82	112.65	77.49	109.74	71.45	69.14	62.08
	MSE	1.43E-1	3.13E-2	5.55E-3	4.26E-5	5.76E-5	1.72E-5	2.57E-3	1.29E-5
	SSIM	4.03E-3	7.46E-1	9.52E-1	9.96E-1	9.98E-1	9.96E-1	9.74E-1	9.96E-1
HEUN	FID	177.35	112.06	112.65	113.39	109.74	102.50	100.08	88.76
	MSE	1.43E-1	7.98E-4	5.55E-3	3.20E-5	5.76E-5	1.13E-5	2.56E-3	5.40E-6
	SSIM	4.03E-3	9.94E-1	9.52E-1	9.99E-1	9.98E-1	9.99E-1	9.77E-1	9.98E-1
LMSD	FID	177.35	112.06	112.65	113.39	109.74	102.50	100.08	88.76
	MSE	1.43E-1	7.98E-4	5.55E-3	3.20E-5	5.76E-5	1.13E-5	2.56E-3	5.40E-6
	SSIM	4.03E-3	9.94E-1	9.52E-1	9.99E-1	9.98E-1	9.99E-1	9.77E-1	9.98E-1

Table 23: CIFAR10 Dataset with Trigger: Stop Sign, target: Corner, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	14.31	8.54	7.80	8.49	8.08	8.17	7.89	7.91
	MSE	7.93E-2	7.56E-2	7.48E-2	7.54E-2	7.32E-2	6.97E-2	6.83E-2	6.21E-2
	SSIM	7.10E-2	8.99E-2	1.03E-1	8.87E-2	9.95E-2	1.13E-1	1.08E-1	1.44E-1
UNIPC	FID	185.20	121.16	119.11	121.15	118.63	118.40	124.84	124.46
	MSE	1.11E-1	1.19E-2	1.07E-2	1.89E-3	2.14E-2	2.65E-2	2.36E-3	4.50E-3
	SSIM	8.50E-3	6.40E-1	6.97E-1	8.36E-1	5.53E-1	5.21E-1	8.73E-1	8.40E-1
DPM. O2	FID	185.20	121.16	119.11	121.15	118.63	118.40	124.84	124.46
	MSE	1.11E-1	1.19E-2	1.07E-2	1.89E-3	2.14E-2	2.65E-2	2.36E-3	4.50E-3
	SSIM	8.50E-3	6.40E-1	6.97E-1	8.36E-1	5.53E-1	5.21E-1	8.73E-1	8.40E-1
DPM. O3	FID	185.20	121.16	119.11	121.15	118.63	118.40	124.84	124.46
	MSE	1.11E-1	1.19E-2	1.07E-2	1.89E-3	2.14E-2	2.65E-2	2.36E-3	4.50E-3
	SSIM	8.50E-3	6.40E-1	6.97E-1	8.36E-1	5.53E-1	5.21E-1	8.73E-1	8.40E-1
DPM++. O2	FID	185.20	121.16	119.11	121.15	118.63	118.40	124.84	124.46
	MSE	1.11E-1	1.19E-2	1.07E-2	1.89E-3	2.14E-2	2.65E-2	2.36E-3	4.50E-3
	SSIM	8.50E-3	6.40E-1	6.97E-1	8.36E-1	5.53E-1	5.21E-1	8.73E-1	8.40E-1
DPM++. O3	FID	185.20	121.16	119.11	121.15	118.63	82.22	124.84	124.46
	MSE	1.11E-1	1.19E-2	1.07E-2	1.89E-3	2.14E-2	2.42E-2	2.36E-3	4.50E-3
	SSIM	8.50E-3	6.40E-1	6.97E-1	8.36E-1	5.53E-1	5.78E-1	8.73E-1	8.40E-1
DEIS	FID	185.20	121.16	119.11	121.15	118.63	118.40	124.84	124.46
	MSE	1.11E-1	1.19E-2	1.07E-2	1.89E-3	2.14E-2	2.65E-2	2.36E-3	4.50E-3
	SSIM	8.50E-3	6.40E-1	6.97E-1	8.36E-1	5.53E-1	5.21E-1	8.73E-1	8.40E-1
PNM	FID	177.35	86.72	124.80	127.35	122.77	122.23	89.46	89.20
	MSE	1.10E-1	5.75E-2	2.91E-2	3.73E-3	4.75E-2	4.18E-2	1.01E-3	2.17E-4
	SSIM	6.74E-3	4.26E-1	7.13E-1	9.62E-1	5.39E-1	5.68E-1	9.78E-1	9.85E-1
HEUN	FID	177.35	129.05	124.80	127.35	122.77	122.23	129.63	133.33
	MSE	1.10E-1	4.46E-2	2.91E-2	3.73E-3	4.75E-2	4.18E-2	1.08E-3	2.11E-4
	SSIM	6.74E-3	5.72E-1	7.13E-1	9.62E-1	5.39E-1	5.68E-1	9.83E-1	9.91E-1
LMSD	FID	177.35	129.05	124.80	127.35	122.77	122.23	129.63	133.33
	MSE	1.10E-1	4.46E-2	2.91E-2	3.73E-3	4.75E-2	4.18E-2	1.08E-3	2.11E-4
	SSIM	6.74E-3	5.72E-1	7.13E-1	9.62E-1	5.39E-1	5.68E-1	9.83E-1	9.91E-1

Table 24: CIFAR10 Dataset with Trigger: Stop Sign, target: Shoe, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
UNIPC	FID	185.20	107.88	108.12	109.60	105.08	109.91	104.87	98.41
	MSE	3.20E-1	3.56E-2	6.73E-3	1.27E-3	5.94E-3	7.13E-4	1.28E-3	1.94E-3
	SSIM	6.07E-3	8.50E-1	9.50E-1	9.72E-1	9.60E-1	9.83E-1	9.77E-1	9.73E-1
DPM. O2	FID	185.20	107.88	108.12	109.60	105.08	109.91	104.87	98.41
	MSE	3.20E-1	3.56E-2	6.73E-3	1.27E-3	5.94E-3	7.13E-4	1.28E-3	1.94E-3
	SSIM	6.07E-3	8.50E-1	9.50E-1	9.72E-1	9.60E-1	9.83E-1	9.77E-1	9.73E-1
DPM. O3	FID	185.20	107.88	108.12	109.60	105.08	109.91	73.51	98.41
	MSE	3.20E-1	3.56E-2	6.73E-3	1.27E-3	5.94E-3	7.13E-4	8.89E-4	1.94E-3
	SSIM	6.07E-3	8.50E-1	9.50E-1	9.72E-1	9.60E-1	9.83E-1	9.81E-1	9.73E-1
DPM++. O2	FID	185.20	107.88	108.12	109.60	105.08	109.91	104.87	98.41
	MSE	3.20E-1	3.56E-2	6.73E-3	1.27E-3	5.94E-3	7.13E-4	1.28E-3	1.94E-3
	SSIM	6.07E-3	8.50E-1	9.50E-1	9.72E-1	9.60E-1	9.83E-1	9.77E-1	9.73E-1
DPM++. O3	FID	185.20	107.88	108.12	109.60	105.08	109.91	104.87	98.41
	MSE	3.20E-1	3.56E-2	6.73E-3	1.27E-3	5.94E-3	7.13E-4	1.28E-3	1.94E-3
	SSIM	6.07E-3	8.50E-1	9.50E-1	9.72E-1	9.60E-1	9.83E-1	9.77E-1	9.73E-1
DEIS	FID	185.20	107.88	108.12	109.60	105.08	109.91	104.87	98.41
	MSE	3.20E-1	3.56E-2	6.73E-3	1.27E-3	5.94E-3	7.13E-4	1.28E-3	1.94E-3
	SSIM	6.07E-3	8.50E-1	9.50E-1	9.72E-1	9.60E-1	9.83E-1	9.77E-1	9.73E-1
PNDM	FID	177.35	74.67	107.85	73.88	71.59	72.62	105.62	67.34
	MSE	3.24E-1	1.05E-1	1.97E-2	1.70E-3	1.17E-2	3.98E-5	1.53E-5	1.59E-5
	SSIM	4.67E-3	6.26E-1	9.32E-1	9.91E-1	9.53E-1	9.97E-1	9.99E-1	9.97E-1
HEUN	FID	177.35	110.54	107.85	108.84	103.94	107.88	105.62	96.78
	MSE	3.24E-1	8.27E-2	1.97E-2	1.86E-3	1.31E-2	4.19E-5	1.53E-5	9.93E-6
	SSIM	4.67E-3	7.25E-1	9.32E-1	9.92E-1	9.49E-1	9.99E-1	9.99E-1	9.99E-1
LMSD	FID	177.35	110.54	107.85	108.84	103.94	107.88	105.62	96.78
	MSE	3.24E-1	8.27E-2	1.97E-2	1.86E-3	1.31E-2	4.19E-5	1.53E-5	9.93E-6
	SSIM	4.67E-3	7.25E-1	9.32E-1	9.92E-1	9.49E-1	9.99E-1	9.99E-1	9.99E-1



Table 25: CIFAR10 Dataset with Trigger: Stop Sign, target: Hat, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	14.31	8.31	7.53	8.10	7.64	7.63	7.63	7.71
	MSE	1.76E-1	1.67E-1	1.66E-1	1.68E-1	1.67E-1	1.62E-1	1.58E-1	1.48E-1
	SSIM	3.41E-2	4.26E-2	4.36E-2	4.05E-2	4.37E-2	4.70E-2	4.72E-2	5.16E-2
UNIPC	FID	185.20	110.69	109.15	110.76	107.63	110.05	106.94	100.60
	MSE	2.33E-1	8.32E-3	6.79E-3	5.52E-3	6.54E-3	1.88E-3	3.30E-3	5.54E-3
	SSIM	8.04E-3	9.06E-1	9.45E-1	9.61E-1	9.55E-1	9.84E-1	9.80E-1	9.73E-1
DPM. O2	FID	185.20	110.69	109.15	110.76	107.63	110.05	106.94	100.60
	MSE	2.33E-1	8.32E-3	6.79E-3	5.52E-3	6.54E-3	1.88E-3	3.30E-3	5.54E-3
	SSIM	8.04E-3	9.06E-1	9.45E-1	9.61E-1	9.55E-1	9.84E-1	9.80E-1	9.73E-1
DPM. O3	FID	185.20	110.69	109.15	110.76	107.63	110.05	106.94	100.60
	MSE	2.33E-1	8.32E-3	6.79E-3	5.52E-3	6.54E-3	1.88E-3	3.30E-3	5.54E-3
	SSIM	8.04E-3	9.06E-1	9.45E-1	9.61E-1	9.55E-1	9.84E-1	9.80E-1	9.73E-1
DPM++. O2	FID	185.20	110.69	109.15	110.76	107.63	110.05	106.94	100.60
	MSE	2.33E-1	8.32E-3	6.79E-3	5.52E-3	6.54E-3	1.88E-3	3.30E-3	5.54E-3
	SSIM	8.04E-3	9.06E-1	9.45E-1	9.61E-1	9.55E-1	9.84E-1	9.80E-1	9.73E-1
DPM++. O3	FID	185.20	110.69	109.15	110.76	107.63	110.05	106.94	100.60
	MSE	2.33E-1	8.32E-3	6.79E-3	5.52E-3	6.54E-3	1.88E-3	3.30E-3	5.54E-3
	SSIM	8.04E-3	9.06E-1	9.45E-1	9.61E-1	9.55E-1	9.84E-1	9.80E-1	9.73E-1
DEIS	FID	185.20	110.69	109.15	110.76	107.63	110.05	106.94	100.60
	MSE	2.33E-1	8.32E-3	6.79E-3	5.52E-3	6.54E-3	1.88E-3	3.30E-3	5.54E-3
	SSIM	8.04E-3	9.06E-1	9.45E-1	9.61E-1	9.55E-1	9.84E-1	9.80E-1	9.73E-1
PNM	FID	177.35	77.16	42.34	75.69	42.47	106.99	8.11	95.57
	MSE	2.34E-1	9.55E-3	3.38E-2	3.55E-3	1.37E-4	8.03E-4	3.04E-5	6.80E-6
	SSIM	7.73E-3	9.49E-1	8.26E-1	9.82E-1	9.98E-1	9.95E-1	9.98E-1	1.00E+0
HEUN	FID	177.35	111.49	110.30	110.51	108.51	106.99	99.88	95.57
	MSE	2.34E-1	7.92E-3	1.04E-2	4.58E-3	9.40E-3	8.03E-4	1.18E-5	6.80E-6
	SSIM	7.73E-3	9.63E-1	9.52E-1	9.77E-1	9.54E-1	9.95E-1	1.00E+0	1.00E+0
LMSD	FID	177.35	111.49	110.30	110.51	108.51	106.99	99.88	95.57
	MSE	2.34E-1	7.92E-3	1.04E-2	4.58E-3	9.40E-3	8.03E-4	1.18E-5	6.80E-6
	SSIM	7.73E-3	9.63E-1	9.52E-1	9.77E-1	9.54E-1	9.95E-1	1.00E+0	1.00E+0

Table 26: CIFAR10 Dataset with Trigger: Grey Box, target: No Shift, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
UNIPC	FID	185.20	121.50	121.93	130.83	123.83	134.04	146.21	159.08
	MSE	1.20E-1	1.18E-3	4.93E-4	4.28E-4	1.05E-3	1.23E-3	1.64E-3	3.83E-3
	SSIM	8.51E-4	9.68E-1	9.76E-1	9.78E-1	9.62E-1	9.62E-1	9.62E-1	9.31E-1
DPM. O2	FID	185.20	121.50	121.93	130.83	123.83	134.04	146.21	159.08
	MSE	1.20E-1	1.18E-3	4.93E-4	4.28E-4	1.05E-3	1.23E-3	1.64E-3	3.83E-3
	SSIM	8.51E-4	9.68E-1	9.76E-1	9.78E-1	9.62E-1	9.62E-1	9.62E-1	9.31E-1
DPM. O3	FID	185.20	121.50	84.27	130.83	123.83	134.04	146.21	159.08
	MSE	1.20E-1	1.18E-3	1.13E-3	4.28E-4	1.05E-3	1.23E-3	1.64E-3	3.83E-3
	SSIM	8.51E-4	9.68E-1	9.55E-1	9.78E-1	9.62E-1	9.62E-1	9.62E-1	9.31E-1
DPM++. O2	FID	185.20	121.50	121.93	130.83	123.83	134.04	146.21	159.08
	MSE	1.20E-1	1.18E-3	4.93E-4	4.28E-4	1.05E-3	1.23E-3	1.64E-3	3.83E-3
	SSIM	8.51E-4	9.68E-1	9.76E-1	9.78E-1	9.62E-1	9.62E-1	9.62E-1	9.31E-1
DPM++. O3	FID	185.20	121.50	121.93	130.83	123.83	134.04	146.21	159.08
	MSE	1.20E-1	1.18E-3	4.93E-4	4.28E-4	1.05E-3	1.23E-3	1.64E-3	3.83E-3
	SSIM	8.51E-4	9.68E-1	9.76E-1	9.78E-1	9.62E-1	9.62E-1	9.62E-1	9.31E-1
DEIS	FID	185.20	121.50	121.93	130.83	123.83	134.04	146.21	159.08
	MSE	1.20E-1	1.18E-3	4.93E-4	4.28E-4	1.05E-3	1.23E-3	1.64E-3	3.83E-3
	SSIM	8.51E-4	9.68E-1	9.76E-1	9.78E-1	9.62E-1	9.62E-1	9.62E-1	9.31E-1
PNDM	FID	177.35	84.49	8.34	54.96	91.59	98.41	154.06	75.09
	MSE	1.20E-1	1.80E-2	1.37E-2	2.36E-4	4.39E-5	1.61E-5	7.00E-6	2.16E-5
	SSIM	9.50E-4	7.88E-1	7.93E-1	9.80E-1	9.90E-1	9.91E-1	9.96E-1	9.85E-1
HEUN	FID	177.35	126.85	127.51	137.93	129.76	143.02	154.06	162.09
	MSE	1.20E-1	2.98E-3	1.27E-3	6.22E-5	3.67E-5	8.34E-6	7.00E-6	7.66E-6
	SSIM	9.50E-4	9.72E-1	9.85E-1	9.96E-1	9.96E-1	9.96E-1	9.96E-1	9.95E-1
LMSD	FID	177.35	126.85	127.51	137.93	129.76	143.02	154.06	162.09
	MSE	1.20E-1	2.98E-3	1.27E-3	6.22E-5	3.67E-5	8.34E-6	7.00E-6	7.66E-6
	SSIM	9.50E-4	9.72E-1	9.85E-1	9.96E-1	9.96E-1	9.96E-1	9.96E-1	9.95E-1

Table 27: CIFAR10 Dataset with Trigger: Grey Box, target: Shift, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
UNIPC	FID	185.20	123.67	117.55	118.87	119.92	122.32	126.07	131.79
	MSE	1.20E-1	2.83E-3	1.61E-3	4.57E-4	2.76E-4	1.82E-3	2.07E-3	4.21E-3
	SSIM	7.09E-4	9.50E-1	9.60E-1	9.76E-1	9.84E-1	9.54E-1	9.55E-1	9.27E-1
DPM. O2	FID	185.20	123.67	117.55	118.87	119.92	122.32	126.07	131.79
	MSE	1.20E-1	2.83E-3	1.61E-3	4.57E-4	2.76E-4	1.82E-3	2.07E-3	4.21E-3
	SSIM	7.09E-4	9.50E-1	9.60E-1	9.76E-1	9.84E-1	9.54E-1	9.55E-1	9.27E-1
DPM. O3	FID	185.20	123.67	117.55	118.87	119.92	122.32	126.07	131.79
	MSE	1.20E-1	2.83E-3	1.61E-3	4.57E-4	2.76E-4	1.82E-3	2.07E-3	4.21E-3
	SSIM	7.09E-4	9.50E-1	9.60E-1	9.76E-1	9.84E-1	9.54E-1	9.55E-1	9.27E-1
DPM++. O2	FID	185.20	123.67	117.55	118.87	119.92	122.32	126.07	131.79
	MSE	1.20E-1	2.83E-3	1.61E-3	4.57E-4	2.76E-4	1.82E-3	2.07E-3	4.21E-3
	SSIM	7.09E-4	9.50E-1	9.60E-1	9.76E-1	9.84E-1	9.54E-1	9.55E-1	9.27E-1
DPM++. O3	FID	185.20	123.67	80.34	118.87	119.92	122.32	126.07	131.79
	MSE	1.20E-1	2.83E-3	1.07E-2	4.57E-4	2.76E-4	1.82E-3	2.07E-3	4.21E-3
	SSIM	7.09E-4	9.50E-1	8.34E-1	9.76E-1	9.84E-1	9.54E-1	9.55E-1	9.27E-1
DEIS	FID	185.20	123.67	117.55	118.87	119.92	122.32	126.07	131.79
	MSE	1.20E-1	2.83E-3	1.61E-3	4.57E-4	2.76E-4	1.82E-3	2.07E-3	4.21E-3
	SSIM	7.09E-4	9.50E-1	9.60E-1	9.76E-1	9.84E-1	9.54E-1	9.55E-1	9.27E-1
PNDM	FID	177.35	88.59	81.09	47.93	86.19	88.97	129.77	102.72
	MSE	1.20E-1	1.71E-2	3.19E-3	2.15E-4	3.71E-4	1.64E-5	6.98E-6	1.42E-5
	SSIM	6.17E-4	8.09E-1	9.57E-1	9.86E-1	9.90E-1	9.94E-1	9.97E-1	9.93E-1
HEUN	FID	177.35	125.94	119.25	121.26	122.38	125.11	129.77	138.51
	MSE	1.20E-1	8.17E-3	2.87E-3	2.09E-4	3.83E-4	7.62E-6	6.98E-6	7.41E-6
	SSIM	6.17E-4	9.27E-1	9.73E-1	9.96E-1	9.94E-1	9.98E-1	9.97E-1	9.97E-1
LMSD	FID	177.35	125.94	119.25	121.26	122.38	125.11	129.77	138.51
	MSE	1.20E-1	8.17E-3	2.87E-3	2.09E-4	3.83E-4	7.62E-6	6.98E-6	7.41E-6
	SSIM	6.17E-4	9.27E-1	9.73E-1	9.96E-1	9.94E-1	9.98E-1	9.97E-1	9.97E-1

Table 28: CIFAR10 Dataset with Trigger: Grey Box, target: Corner, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	14.31	9.91	10.94	12.99	15.06	19.85	28.11	53.35
	MSE	7.86E-2	5.56E-2	5.34E-2	4.97E-2	5.01E-2	3.87E-2	2.74E-2	1.32E-2
	SSIM	7.17E-2	2.50E-1	2.80E-1	3.29E-1	3.35E-1	4.60E-1	5.88E-1	7.73E-1
UNIPC	FID	185.20	125.95	127.27	124.84	125.84	131.79	134.52	136.83
	MSE	1.05E-1	2.18E-3	8.14E-4	5.05E-4	4.90E-4	1.09E-3	1.69E-3	5.51E-3
	SSIM	1.13E-3	9.16E-1	9.33E-1	9.48E-1	9.54E-1	9.36E-1	9.24E-1	8.36E-1
DPM. O2	FID	185.20	125.95	127.27	124.84	125.84	131.79	134.52	136.83
	MSE	1.05E-1	2.18E-3	8.14E-4	5.05E-4	4.90E-4	1.09E-3	1.69E-3	5.51E-3
	SSIM	1.13E-3	9.16E-1	9.33E-1	9.48E-1	9.54E-1	9.36E-1	9.24E-1	8.36E-1
DPM. O3	FID	185.20	125.95	127.27	124.84	125.84	131.79	134.52	136.83
	MSE	1.05E-1	2.18E-3	8.14E-4	5.05E-4	4.90E-4	1.09E-3	1.69E-3	5.51E-3
	SSIM	1.13E-3	9.16E-1	9.33E-1	9.48E-1	9.54E-1	9.36E-1	9.24E-1	8.36E-1
DPM++. O2	FID	185.20	125.95	127.27	124.84	125.84	131.79	134.52	136.83
	MSE	1.05E-1	2.18E-3	8.14E-4	5.05E-4	4.90E-4	1.09E-3	1.69E-3	5.51E-3
	SSIM	1.13E-3	9.16E-1	9.33E-1	9.48E-1	9.54E-1	9.36E-1	9.24E-1	8.36E-1
DPM++. O3	FID	185.20	125.95	127.27	124.84	125.84	131.79	134.52	136.83
	MSE	1.05E-1	2.18E-3	8.14E-4	5.05E-4	4.90E-4	1.09E-3	1.69E-3	5.51E-3
	SSIM	1.13E-3	9.16E-1	9.33E-1	9.48E-1	9.54E-1	9.36E-1	9.24E-1	8.36E-1
DEIS	FID	185.20	125.95	127.27	124.84	125.84	131.79	134.52	136.83
	MSE	1.05E-1	2.18E-3	8.14E-4	5.05E-4	4.90E-4	1.09E-3	1.69E-3	5.51E-3
	SSIM	1.13E-3	9.16E-1	9.33E-1	9.48E-1	9.54E-1	9.36E-1	9.24E-1	8.36E-1
PNM	FID	177.35	134.07	93.84	49.66	94.22	102.26	147.95	156.45
	MSE	1.05E-1	3.56E-3	4.82E-3	9.19E-5	3.88E-5	2.01E-5	9.18E-6	7.25E-6
	SSIM	1.38E-3	9.66E-1	8.97E-1	9.78E-1	9.87E-1	9.88E-1	9.94E-1	9.94E-1
HEUN	FID	177.35	134.07	135.99	134.72	137.73	145.55	147.95	156.45
	MSE	1.05E-1	3.56E-3	4.92E-4	4.11E-5	5.36E-5	1.44E-5	9.18E-6	7.25E-6
	SSIM	1.38E-3	9.66E-1	9.92E-1	9.96E-1	9.95E-1	9.94E-1	9.94E-1	9.94E-1
LMSD	FID	177.35	134.07	135.99	134.72	137.73	145.55	147.95	156.45
	MSE	1.05E-1	3.56E-3	4.92E-4	4.11E-5	5.36E-5	1.44E-5	9.18E-6	7.25E-6
	SSIM	1.38E-3	9.66E-1	9.92E-1	9.96E-1	9.95E-1	9.94E-1	9.94E-1	9.94E-1

Table 29: CIFAR10 Dataset with Trigger: Grey Box, target: Shoe, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
UNIPC	FID	185.20	108.96	109.27	113.59	110.68	111.37	105.23	105.27
	MSE	3.37E-1	2.20E-3	1.02E-3	8.76E-4	9.58E-4	1.80E-3	1.97E-3	3.59E-3
	SSIM	2.29E-4	9.80E-1	9.82E-1	9.82E-1	9.83E-1	9.76E-1	9.76E-1	9.65E-1
DPM. O2	FID	185.20	108.96	109.27	113.59	110.68	111.37	105.23	105.27
	MSE	3.37E-1	2.20E-3	1.02E-3	8.76E-4	9.58E-4	1.80E-3	1.97E-3	3.59E-3
	SSIM	2.29E-4	9.80E-1	9.82E-1	9.82E-1	9.83E-1	9.76E-1	9.76E-1	9.65E-1
DPM. O3	FID	185.20	108.96	109.27	113.59	110.68	111.37	105.23	105.27
	MSE	3.37E-1	2.20E-3	1.02E-3	8.76E-4	9.58E-4	1.80E-3	1.97E-3	3.59E-3
	SSIM	2.29E-4	9.80E-1	9.82E-1	9.82E-1	9.83E-1	9.76E-1	9.76E-1	9.65E-1
DPM++. O2	FID	185.20	75.06	109.27	113.59	110.68	111.37	105.23	105.27
	MSE	3.37E-1	1.96E-2	1.02E-3	8.76E-4	9.58E-4	1.80E-3	1.97E-3	3.59E-3
	SSIM	2.29E-4	8.92E-1	9.82E-1	9.82E-1	9.83E-1	9.76E-1	9.76E-1	9.65E-1
DPM++. O3	FID	185.20	108.96	109.27	113.59	110.68	111.37	105.23	105.27
	MSE	3.37E-1	2.20E-3	1.02E-3	8.76E-4	9.58E-4	1.80E-3	1.97E-3	3.59E-3
	SSIM	2.29E-4	9.80E-1	9.82E-1	9.82E-1	9.83E-1	9.76E-1	9.76E-1	9.65E-1
DEIS	FID	185.20	108.96	109.27	113.59	110.68	111.37	105.23	105.27
	MSE	3.37E-1	2.20E-3	1.02E-3	8.76E-4	9.58E-4	1.80E-3	1.97E-3	3.59E-3
	SSIM	2.29E-4	9.80E-1	9.82E-1	9.82E-1	9.83E-1	9.76E-1	9.76E-1	9.65E-1
PNDM	FID	177.35	76.12	77.07	44.45	74.86	76.19	73.59	73.31
	MSE	3.37E-1	6.60E-2	4.31E-3	8.33E-5	5.15E-5	2.28E-5	1.95E-5	1.49E-5
	SSIM	2.08E-4	7.43E-1	9.75E-1	9.95E-1	9.97E-1	9.97E-1	9.98E-1	9.97E-1
HEUN	FID	177.35	109.16	111.45	112.72	109.41	108.98	104.22	105.28
	MSE	3.37E-1	4.08E-3	8.50E-4	1.14E-4	5.27E-5	1.76E-5	1.19E-5	8.16E-6
	SSIM	2.08E-4	9.87E-1	9.96E-1	9.98E-1	9.99E-1	9.99E-1	9.99E-1	9.99E-1
LMSD	FID	177.35	109.16	111.45	112.72	109.41	108.98	104.22	105.28
	MSE	3.37E-1	4.08E-3	8.50E-4	1.14E-4	5.27E-5	1.76E-5	1.19E-5	8.16E-6
	SSIM	2.08E-4	9.87E-1	9.96E-1	9.98E-1	9.99E-1	9.99E-1	9.99E-1	9.99E-1

Table 30: CIFAR10 Dataset with Trigger: Grey Box, target: Hat, and inference-time clipping.

Sampler	P.R. Metric	0%	5%	10%	20%	30%	50%	70%	90%
ANCESTRAL	FID	14.31	8.42	8.82	8.89	8.97	10.11	11.32	17.82
	MSE	1.74E-1	1.24E-1	1.08E-1	1.09E-1	1.12E-1	1.01E-1	9.63E-2	8.57E-2
	SSIM	3.43E-2	2.08E-1	2.83E-1	2.82E-1	2.66E-1	3.26E-1	3.55E-1	4.07E-1
UNIPC	FID	185.20	110.41	111.65	109.41	112.08	111.42	105.79	107.68
	MSE	2.40E-1	5.12E-4	3.63E-4	5.34E-4	8.85E-4	1.63E-3	3.46E-3	6.46E-3
	SSIM	2.97E-4	9.88E-1	9.92E-1	9.92E-1	9.92E-1	9.88E-1	9.80E-1	9.70E-1
DPM. O2	FID	185.20	110.41	111.65	109.41	112.08	111.42	105.79	107.68
	MSE	2.40E-1	5.12E-4	3.63E-4	5.34E-4	8.85E-4	1.63E-3	3.46E-3	6.46E-3
	SSIM	2.97E-4	9.88E-1	9.92E-1	9.92E-1	9.92E-1	9.88E-1	9.80E-1	9.70E-1
DPM. O3	FID	185.20	110.41	111.65	109.41	112.08	111.42	105.79	107.68
	MSE	2.40E-1	5.12E-4	3.63E-4	5.34E-4	8.85E-4	1.63E-3	3.46E-3	6.46E-3
	SSIM	2.97E-4	9.88E-1	9.92E-1	9.92E-1	9.92E-1	9.88E-1	9.80E-1	9.70E-1
DPM++. O2	FID	185.20	110.41	111.65	109.41	112.08	111.42	105.79	107.68
	MSE	2.40E-1	5.12E-4	3.63E-4	5.34E-4	8.85E-4	1.63E-3	3.46E-3	6.46E-3
	SSIM	2.97E-4	9.88E-1	9.92E-1	9.92E-1	9.92E-1	9.88E-1	9.80E-1	9.70E-1
DPM++. O3	FID	185.20	110.41	111.65	109.41	112.08	111.42	105.79	107.68
	MSE	2.40E-1	5.12E-4	3.63E-4	5.34E-4	8.85E-4	1.63E-3	3.46E-3	6.46E-3
	SSIM	2.97E-4	9.88E-1	9.92E-1	9.92E-1	9.92E-1	9.88E-1	9.80E-1	9.70E-1
DEIS	FID	185.20	110.41	111.65	109.41	112.08	111.42	105.79	107.68
	MSE	2.40E-1	5.12E-4	3.63E-4	5.34E-4	8.85E-4	1.63E-3	3.46E-3	6.46E-3
	SSIM	2.97E-4	9.88E-1	9.92E-1	9.92E-1	9.92E-1	9.88E-1	9.80E-1	9.70E-1
PNM	FID	177.35	112.35	111.44	77.27	41.78	77.43	101.54	102.53
	MSE	2.40E-1	4.60E-4	9.89E-5	7.03E-5	3.09E-5	2.03E-5	7.09E-6	6.13E-6
	SSIM	1.79E-4	9.97E-1	9.99E-1	9.98E-1	9.98E-1	9.99E-1	1.00E+0	1.00E+0
HEUN	FID	177.35	112.35	111.44	109.32	110.74	109.52	101.54	102.53
	MSE	2.40E-1	4.60E-4	9.89E-5	2.39E-5	1.41E-5	1.33E-5	7.09E-6	6.13E-6
	SSIM	1.79E-4	9.97E-1	9.99E-1	9.99E-1	9.99E-1	1.00E+0	1.00E+0	1.00E+0
LMSD	FID	177.35	112.35	111.44	109.32	110.74	109.52	101.54	102.53
	MSE	2.40E-1	4.60E-4	9.89E-5	2.39E-5	1.41E-5	1.33E-5	7.09E-6	6.13E-6
	SSIM	1.79E-4	9.97E-1	9.99E-1	9.99E-1	9.99E-1	1.00E+0	1.00E+0	1.00E+0

Table 31: DDPM performs on **Blur**, **Line**, **Box**, and **Box** with CIFAR10 Dataset, trigger: Stop Sign, and target: Shoe.

Sampler	P.R. Metric	0%	10%	20%	30%	50%	70%	90%
UNIPC, Blur	LPIPS	3.25E-1	2.78E-1	2.77E-1	2.68E-1	2.56E-1	2.55E-1	2.59E-1
	MSE	2.97E-1	3.58E-3	1.31E-3	1.61E-3	4.90E-4	1.99E-4	8.39E-5
	SSIM	5.66E-2	9.85E-1	9.94E-1	9.93E-1	9.97E-1	9.99E-1	9.99E-1
UNIPC, Line	LPIPS	3.19E-1	2.78E-1	2.82E-1	2.76E-1	2.66E-1	2.65E-1	2.58E-1
	MSE	3.03E-1	3.78E-3	2.83E-3	3.29E-3	1.27E-3	3.37E-4	6.92E-5
	SSIM	4.61E-2	9.85E-1	9.89E-1	9.87E-1	9.95E-1	9.98E-1	9.99E-1
UNIPC, Box	LPIPS	3.10E-1	2.87E-1	2.85E-1	2.85E-1	2.70E-1	2.63E-1	2.63E-1
	MSE	3.34E-1	2.15E-2	3.17E-2	3.47E-2	3.29E-2	1.93E-2	1.11E-2
	SSIM	1.05E-2	9.10E-1	8.76E-1	8.56E-1	8.67E-1	9.19E-1	9.51E-1
DPM. O2, Blur	LPIPS	3.25E-1	2.79E-1	2.70E-1	2.65E-1	2.57E-1	2.55E-1	2.57E-1
	MSE	2.97E-1	2.98E-3	1.52E-3	1.82E-3	5.58E-4	2.93E-4	1.37E-4
	SSIM	5.66E-2	9.87E-1	9.93E-1	9.91E-1	9.97E-1	9.98E-1	9.99E-1
DPM. O2, Line	LPIPS	3.19E-1	2.88E-1	2.80E-1	2.73E-1	2.67E-1	2.63E-1	2.60E-1
	MSE	3.03E-1	3.02E-3	3.25E-3	3.19E-3	1.39E-3	3.21E-4	1.97E-4
	SSIM	4.61E-2	9.87E-1	9.87E-1	9.88E-1	9.94E-1	9.98E-1	9.99E-1
DPM. O2, Box	LPIPS	3.10E-1	2.91E-1	2.87E-1	2.83E-1	2.71E-1	2.59E-1	2.67E-1
	MSE	3.34E-1	2.07E-2	3.24E-2	3.45E-2	3.60E-2	2.14E-2	1.04E-2
	SSIM	1.05E-2	9.14E-1	8.74E-1	8.57E-1	8.53E-1	9.11E-1	9.50E-1
DPM++. O2, Blur	LPIPS	3.25E-1	2.78E-1	2.76E-1	2.68E-1	2.59E-1	2.52E-1	2.54E-1
	MSE	2.97E-1	3.56E-3	1.52E-3	1.85E-3	4.64E-4	5.16E-4	4.34E-6
	SSIM	5.66E-2	9.84E-1	9.93E-1	9.92E-1	9.98E-1	9.97E-1	1.00E+0
DPM++. O2, Line	LPIPS	3.19E-1	2.87E-1	2.76E-1	2.73E-1	2.67E-1	2.61E-1	2.57E-1
	MSE	3.03E-1	3.55E-3	4.38E-3	4.30E-3	1.68E-3	4.40E-4	9.60E-5
	SSIM	4.61E-2	9.85E-1	9.83E-1	9.83E-1	9.93E-1	9.98E-1	9.99E-1
DPM++. O2, Box	LPIPS	3.10E-1	2.87E-1	2.85E-1	2.79E-1	2.74E-1	2.65E-1	2.64E-1
	MSE	3.34E-1	2.44E-2	3.40E-2	3.12E-2	3.79E-2	1.97E-2	1.08E-2
	SSIM	1.05E-2	8.98E-1	8.72E-1	8.66E-1	8.45E-1	9.15E-1	9.49E-1
DEIS, Blur	LPIPS	3.25E-1	2.81E-1	2.76E-1	2.72E-1	2.60E-1	2.51E-1	2.50E-1
	MSE	2.97E-1	2.87E-3	1.34E-3	2.96E-3	2.02E-4	3.58E-4	1.61E-4
	SSIM	5.66E-2	9.87E-1	9.94E-1	9.88E-1	9.98E-1	9.98E-1	9.99E-1
DEIS, Line	LPIPS	3.19E-1	2.82E-1	2.77E-1	2.76E-1	2.71E-1	2.58E-1	2.61E-1
	MSE	3.03E-1	3.54E-3	2.50E-3	3.31E-3	1.29E-3	2.57E-4	5.03E-6
	SSIM	4.61E-2	9.85E-1	9.90E-1	9.86E-1	9.95E-1	9.98E-1	1.00E+0
DEIS, Box	LPIPS	3.10E-1	2.84E-1	2.85E-1	2.81E-1	2.70E-1	2.61E-1	2.67E-1
	MSE	3.34E-1	2.26E-2	3.26E-2	3.49E-2	3.46E-2	2.05E-2	1.03E-2
	SSIM	1.05E-2	9.07E-1	8.74E-1	8.52E-1	8.59E-1	9.15E-1	9.54E-1

Table 32: DDPM performs on **Blur**, **Line**, **Box**, and **Box** with CIFAR10 Dataset, trigger: Stop Sign, and target: Hat.

Sampler	P.R. Metric	0%	10%	20%	30%	50%	70%	90%
UNIPC, Blur	LPIPS	3.25E-1	2.72E-1	2.66E-1	2.63E-1	2.53E-1	2.39E-1	2.45E-1
	MSE	2.85E-1	1.25E-2	1.63E-3	2.09E-3	3.59E-4	1.01E-3	6.57E-4
	SSIM	2.98E-2	9.52E-1	9.93E-1	9.91E-1	9.98E-1	9.95E-1	9.96E-1
UNIPC, Line	LPIPS	3.19E-1	2.73E-1	2.69E-1	2.65E-1	2.54E-1	2.41E-1	2.45E-1
	MSE	2.83E-1	1.56E-2	2.33E-3	2.71E-3	1.11E-3	8.63E-4	1.23E-3
	SSIM	2.13E-2	9.42E-1	9.90E-1	9.89E-1	9.95E-1	9.96E-1	9.94E-1
UNIPC, Box	LPIPS	3.10E-1	2.87E-1	2.79E-1	2.79E-1	2.65E-1	2.51E-1	2.52E-1
	MSE	3.04E-1	3.63E-2	2.16E-2	1.43E-2	4.97E-3	3.31E-3	6.02E-3
	SSIM	-1.37E-3	8.76E-1	9.23E-1	9.47E-1	9.81E-1	9.87E-1	9.76E-1
DPM. O2, Blur	LPIPS	3.25E-1	2.73E-1	2.71E-1	2.60E-1	2.55E-1	2.41E-1	2.37E-1
	MSE	2.85E-1	1.16E-2	1.52E-3	1.58E-3	8.13E-4	1.07E-3	7.07E-4
	SSIM	2.98E-2	9.58E-1	9.93E-1	9.93E-1	9.96E-1	9.95E-1	9.97E-1
DPM. O2, Line	LPIPS	3.19E-1	2.76E-1	2.68E-1	2.63E-1	2.57E-1	2.42E-1	2.45E-1
	MSE	2.83E-1	1.43E-2	2.28E-3	1.84E-3	1.03E-3	1.13E-3	2.68E-3
	SSIM	2.13E-2	9.46E-1	9.91E-1	9.92E-1	9.95E-1	9.94E-1	9.89E-1
DPM. O2, Box	LPIPS	3.10E-1	2.86E-1	2.76E-1	2.77E-1	2.65E-1	2.54E-1	2.52E-1
	MSE	3.04E-1	3.83E-2	2.15E-2	1.40E-2	3.26E-3	3.50E-3	1.41E-3
	SSIM	-1.37E-3	8.67E-1	9.23E-1	9.48E-1	9.87E-1	9.86E-1	9.93E-1
DPM++. O2, Blur	LPIPS	3.25E-1	2.71E-1	2.65E-1	2.62E-1	2.59E-1	2.39E-1	2.45E-1
	MSE	2.85E-1	1.16E-2	1.46E-3	1.77E-3	7.03E-4	9.69E-4	4.74E-4
	SSIM	2.98E-2	9.56E-1	9.94E-1	9.92E-1	9.96E-1	9.95E-1	9.98E-1
DPM++. O2, Line	LPIPS	3.19E-1	2.68E-1	2.72E-1	2.62E-1	2.52E-1	2.39E-1	2.41E-1
	MSE	2.83E-1	1.65E-2	1.57E-3	1.78E-3	1.38E-4	8.54E-4	2.20E-3
	SSIM	2.13E-2	9.38E-1	9.92E-1	9.92E-1	9.99E-1	9.96E-1	9.90E-1
DPM++. O2, Box	LPIPS	3.10E-1	2.85E-1	2.79E-1	2.72E-1	2.67E-1	2.52E-1	2.53E-1
	MSE	3.04E-1	4.04E-2	2.17E-2	1.31E-2	4.82E-3	4.04E-3	6.23E-3
	SSIM	-1.37E-3	8.62E-1	9.22E-1	9.52E-1	9.81E-1	9.84E-1	9.73E-1
DEIS, Blur	LPIPS	3.25E-1	2.69E-1	2.71E-1	2.66E-1	2.55E-1	2.43E-1	2.46E-1
	MSE	2.85E-1	1.11E-2	1.20E-3	1.85E-3	9.45E-4	5.66E-4	4.59E-4
	SSIM	2.98E-2	9.59E-1	9.94E-1	9.93E-1	9.95E-1	9.97E-1	9.97E-1
DEIS, Line	LPIPS	3.19E-1	2.73E-1	2.66E-1	2.59E-1	2.53E-1	2.43E-1	2.41E-1
	MSE	2.83E-1	1.39E-2	2.92E-3	2.24E-3	5.60E-4	1.29E-3	1.98E-3
	SSIM	2.13E-2	9.49E-1	9.88E-1	9.91E-1	9.97E-1	9.94E-1	9.90E-1
DEIS, Box	LPIPS	3.10E-1	2.84E-1	2.79E-1	2.72E-1	2.65E-1	2.53E-1	2.44E-1
	MSE	3.04E-1	4.07E-2	2.16E-2	1.16E-2	4.20E-3	3.75E-3	5.07E-3
	SSIM	-1.37E-3	8.62E-1	9.23E-1	9.57E-1	9.83E-1	9.85E-1	9.79E-1