
Invisible Backdoor Attacks on Diffusion Models

Sen Li¹, Junchi Ma², Minhao Cheng³

¹The Hong Kong University of Science and Technology, Hong Kong, China

²The University of British Columbia, Canada

³Penn State University, USA

slien@connect.ust.hk, junchima@student.ubc.ca, mmc7149@psu.edu

Abstract

In recent years, diffusion models have achieved remarkable success in the realm of high-quality image generation, garnering increased attention. This surge in interest is paralleled by a growing concern over the security threats associated with diffusion models, largely attributed to their susceptibility to malicious exploitation. Notably, recent research has brought to light the vulnerability of diffusion models to backdoor attacks, enabling the generation of specific target images through corresponding triggers. However, prevailing backdoor attack methods rely on manually crafted trigger generation functions, often manifesting as discernible patterns incorporated into input noise, thus rendering them susceptible to human detection. In this paper, we present an innovative and versatile optimization framework designed to acquire invisible triggers, enhancing the stealthiness and resilience of inserted backdoors. Our proposed framework is applicable to both unconditional and conditional diffusion models, and notably, we are the pioneers in demonstrating the backdooring of diffusion models within the context of text-guided image editing and inpainting pipelines. Moreover, we also show that the backdoors in the conditional generation can be directly applied to model watermarking for model ownership verification, which further boosts the significance of the proposed framework. Extensive experiments on various commonly used samplers and datasets verify the efficacy and stealthiness of the proposed framework. Our code is publicly available at https://github.com/invisibleTriggerDiffusion/invisible_triggers_for_diffusion.

1 Introduction

Recently, diffusion models have showcased exceptional performance in generating high-quality and diverse image [13, 22, 9, 14, 21]. Based on diffusion models, many popular applications have been developed including GLIDE [21], Imagen [25], and Stable Diffusion [24]. These applications serve as powerful tools for unleashing creativity in content generation. By slowly adding random noise to data, diffusion models learn to reverse the diffusion process to construct desired data samples from the noise. While this approach has proven excellent in creative content generation, it concurrently introduces novel security challenges that warrant careful consideration.

Instances of backdoor attacks on diffusion models, as explored in previous studies [7, 6, 8, 27], have illuminated the potential threats associated with manipulating the diffusion process. However, prevailing approaches either incorporate conspicuous image triggers and additional text into the input noise or prompt to backdoor diffusion models [7, 6, 8]. Alternatively, some methods focus solely on substituting input text characters with non-Latin characters, impacting text encoders rather than diffusion models [27]. While these strategies have demonstrated commendable success rates, the visibility of triggers in [7] and [6] renders them susceptible to detection through human inspection. Notably, none of the prior works, to the best of our knowledge, have ventured into the realm of invisible image triggers for backdooring diffusion models. The adaptability of invisible triggers to

different input conditions allows them to seamlessly blend into the background noise, enhancing their robustness against diverse inputs. Moreover, the subtlety of invisible triggers contributes to a more sustained and persistent backdoor presence. This subtlety facilitates prolonged content manipulation, presenting a challenge for defensive mechanisms to promptly identify and effectively mitigate the threat. The nuanced nature of invisible triggers, therefore, adds an extra layer of complexity and resilience to backdoor attacks on diffusion models.

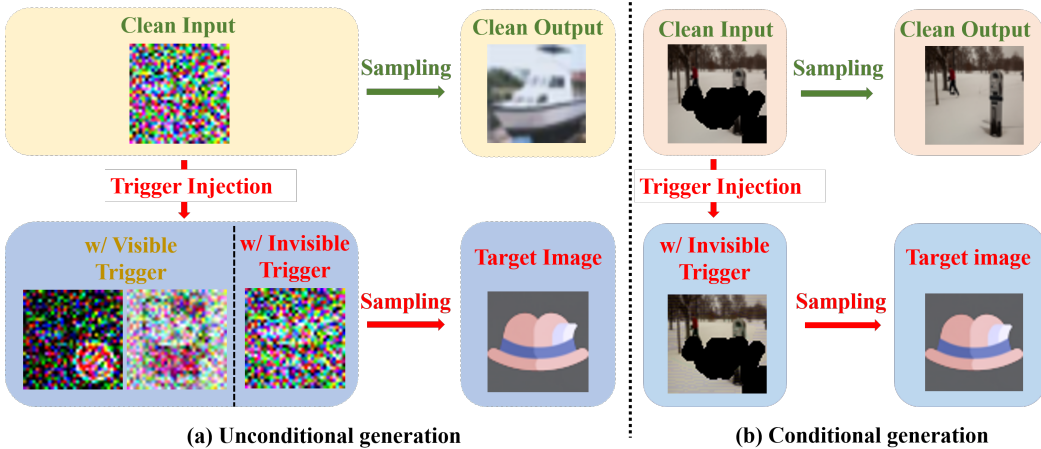


Figure 1: Illustration of our proposed invisible triggers, and visible triggers used in [7, 8, 6].

In this paper, illustrated in Figure 1, we delve into the realm of backdoor attacks featuring invisible image triggers applicable to both unconditional and conditional diffusion models. Our approach introduces a pioneering and comprehensive framework, formulated through bi-level optimization, designed to learn input-aware invisible triggers. The inner optimization phase involves optimizing a trigger generator, given a fixed diffusion model, to enable the generation of target images while maintaining the imperceptibility of the trigger. Simultaneously, the outer optimization phase, given a fixed trigger generator, optimizes the diffusion model to exhibit strong performance on both clean and poisoned data. This framework seamlessly integrates backdooring for both conditional and unconditional diffusion models within a unified optimization problem.

Specifically, when the prior is a random Gaussian noise for the unconditional diffusion model, we train the diffusion model to learn different distribution mappings so that the diffusion model would generate the target image given an initial noise that contains the trigger. Importantly, our framework can insert multiple trigger-target pairs and generate instance-adaptive initial noise rather than applying a uniform perturbation. For conditional diffusion models, we employ a simple neural network to generate instance-adaptive and invisible triggers, effectively incorporating backdoors into additional priors. This enables the conditional diffusion model to generate the target image regardless of the given text, representing a novel contribution compared to previous works focusing on trigger design in the text domain. The invisibility of our proposed backdoor trigger also provides us a seamlessly model watermarking method for model ownership verification. If the suspected model is derived from the watermarked model, once the trigger is being activated, the model would generate a designated target image regardless of any prompt/instruction given. Unlike the previous watermark method using prompt word as trigger, our proposed watermark is robust against different kinds of image editing and instructions.

In summary, our work marks a significant milestone as we introduce a novel and versatile optimization framework to inject input-aware invisible image triggers into both unconditional and conditional diffusion models. This enhancement renders the injected backdoor more covert and robust.

2 Related work

Diffusion models As a new family of powerful generative models, diffusion models could achieve superb performance on high-quality image synthesis [13, 21, 25, 24]. They have shown impressive results on various tasks, such as class-to-image generation [9, 14], text-to-image generation [25], image-to-image translation [20], text-guided image editing/inpainting [21, 24], and so on. A more detailed introduction to diffusion models can be found in Appendix D. With the powerful capability,

the research community has started to focus on the potential security issues that diffusion models may introduce. In this paper, we propose a strong attack framework which can make diffusion models perform maliciously when some invisible patterns are injected into the input, revealing the potential severe security risk that previous works did not cover.

Backdoor attacks on diffusion models Recently, diffusion models have been shown to be vulnerable to backdoor attacks [7, 6, 8, 27]. Struppek et al. [27] proposed to backdoor the text encoder only in text-to-image diffusion models by replacing text characters with non-Latin characters, showing the potential threat in text-to-image generation. However, their method did not consider the backdoor threat on the diffusion models, limiting the practical use of the method. Chou et al. [7] and Chen et al. [6] proposed to backdoor diffusion models in different ways. Their methods focus on backdooring unconditional diffusion models, which may not be applicable to backdoor diffusion models in conditional case. Very recently, Chou et al. [8] proposed a unified framework on backdoor attack for both unconditional diffusion models and text-to-image diffusion models.

To the best of our knowledge, all previous works are working on visible trigger in the diffusion model which could be easily detected and ruled out. In this paper, we propose a novel and general framework to inject invisible triggers into both unconditional and conditional diffusion models. Although there are previous works utilizing bi-level optimization to generate invisible triggers in classification models [11, 10], the context of learning an invisible backdoor trigger in diffusion models significantly differs. We have included an in-depth discussion in Appendix E. Our work further distinguishes itself by formulating a general loss function capable of accommodating various efficient samplers, such as DDIM [26] and DPMSolver [18], as opposed to being limited to a single sampler as in [7].

3 Methodology

3.1 Threat Model

We adopt a similar threat model as prior research [7, 8], where the attack involves two distinct parties. An *attacker* is responsible for injecting backdoors into diffusion models, subsequently releasing the backdoored models. Meanwhile, a *user* downloads these pre-trained models from the web for practical use and have full access to the backdoored models. Additionally, the *user* possesses a subset of clean data for evaluating the models’ performance. Within this model, the *attacker* retains control over the training procedure of the diffusion, encompassing both initial training and fine-tuning. The *attacker* is also granted the capability to modify the training datasets, allowing the addition of supplementary examples. Throughout the training process, the *attacker* endeavors to release backdoored models that exhibit designated behavior (e.g. generating specific image) when the input is injected with the trigger, while behaving normally on inputs without the trigger. Therefore, the *attacker’s* dual objectives include achieving high specificity, ensuring the backdoored models perform maliciously by generating target images when triggered, and maintaining high utility, signifying performance similar to clean models in generating high-quality images consistent with the training dataset distribution. Since current backdoor triggers all make poisoned images visually different from clean images and universal across all the inputs, they can be detected easily by universal perturbation-based detection and human inspection. In this paper, we aim to learn input-aware invisible triggers to make the injected backdoor stronger and more stealthy.

3.2 Optimization framework for learnable invisible trigger

In this paper, our goal is to inject invisible triggers into both unconditional and conditional diffusion models, enhancing the backdoor’s stealthiness and effectiveness. To address this, we formulate the task as a bi-level optimization problem, which learns the invisible trigger to be inserted given different priors. For the purpose of generality, let g be a trigger generator that generates invisible triggers given different priors P , \mathcal{A} be the trigger insertion function which inserts the invisible trigger generated by g into P , and y be the target image that the backdoored diffusion model will generate when the trigger is activated. Let ϵ_θ be the diffusion model which takes the prior P as input to predict the noise, and \mathcal{S} denote the whole sampling process of diffusion models which takes diffusion model ϵ_θ and P as input to generate real data by iterative sampling. We first show the general optimization framework for different priors (i.e., unconditional and conditional generation). Then we will elaborate on the specific parameterization of g and \mathcal{A} for different priors P in detail.

We formulate learning invisible backdoor for diffusion model into a bi-level optimization problem. In the inner optimization, we optimize the trigger generator g given fixed ϵ_θ to generate the target

image \mathbf{y} . The MSE(mean squared error) is used as loss function to train g as:

$$L_{inner}(\epsilon_\theta, P, g(P)) = \left\| \mathcal{S}(\epsilon_\theta, \mathcal{A}(P, g(P))) - \mathbf{y} \right\|^2. \quad (1)$$

At the same time, to ensure the invisibility of the generated trigger, the generated trigger is bounded by ℓ_p norm (ℓ_∞ used in this paper), where we use PGD [19] optimization to force the constraint. To optimize g , all intermediate results during sampling have to be stored to compute the gradient with respect to g . It thus becomes infeasible to sample many steps like the original DDPM [13] sampling. Alternatively, we use DDIM [26] to perform the accelerated sampling process to make the optimization tractable.

For the outer optimization of bi-level optimization, we optimize the diffusion model ϵ_θ to correctly predict the noise for both clean data and poisoned data (i.e., backdoor diffusion process). Let $L_{outer}(\epsilon_\theta, \mathbf{x}_0, P, g(P), \mathbf{y}, t)$ denote the loss function for the outer optimization, where L_{outer} is used to enforce the backdoor trigger's efficacy and model's utility. We defer the detailed L_{outer} design in the following section based on different priors in condition and uncondition case.

Therefore, the whole bi-level optimization framework for input-aware invisible trigger can be formulated as

$$\begin{aligned} \min_{\theta} L_{outer}(\epsilon_\theta, \mathbf{x}_0, P, g_\theta^*(P), \mathbf{y}, t) \\ \text{s.t. } g_\theta^* = \arg \min_g \left\| \mathcal{S}(\epsilon_\theta, \mathcal{A}(P, g(P))) - \mathbf{y} \right\|^2, \|g(P)\|_\infty \leq C \end{aligned} \quad (2)$$

where $\|\cdot\|_\infty$ denotes the ℓ_∞ norm, and C is the norm bound of the trigger. Given the general framework in Equation 2, we now describe the specific parameterization of the loss used under both unconditioned and conditioned scenarios.

3.2.1 Backdooring unconditional diffusion models

When we backdoor unconditional diffusion model, the prior would be regarded as random Gaussian noise $P = \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Specifically, different from visible trigger used in prior works [7, 8, 6] on backdooring unconditional diffusion model, we optimize g as a universal trigger generator for any random noise sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. That is, $g(\epsilon) = \delta, \forall \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The trigger injection function \mathcal{A} then could be defined as $\mathcal{A}(P, g(P)) = \mathcal{A}(\epsilon, \delta) = \delta + \epsilon$. In other words, we aim to make the diffusion model generate a specific image (target image) given any poisoned noise sampled from $\mathcal{N}(\delta, \mathbf{I})$. By generating different δ , our proposed method could insert multiple invisible universal trigger-target pairs simultaneously. Furthermore, to make the trigger more invisible and versatile, instead of keeping δ universal, we optimize g to be a trigger distribution $\mathcal{N}(\delta, \mathbf{I})$ so that any noise draw from the trigger distribution will make the diffusion model generate target image. Formally, we make $g(\epsilon) = \delta'$ and $\mathcal{A}(P, g(P)) = \delta' + \epsilon, \delta' \sim \mathcal{N}(\delta, \mathbf{I}), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Keeping the trigger invisible, our trigger generator is dynamic and sample-specific, which is able to bypass the current universal perturbation-based detection and defense.

Given the above parameterization of g and \mathcal{A} , we define the loss function L_{outer} based on DDIM [26] where we aim to create a secret mapping for poisoned data between target image \mathbf{y} and poisoned distribution $\mathcal{N}(\delta, \mathbf{I})$ or $\mathcal{N}(\delta', \mathbf{I})$. Let $\mathbf{x}'_0 = \mathbf{y}$ (the target image distribution). $\mathbf{x}'_1, \mathbf{x}'_2, \dots, \mathbf{x}'_T \in \mathbb{R}^d$ is then generated by gradually adding noise into \mathbf{x}'_0 where noise schedule is also controlled by β_t . The backdoored forward process is defined as:

$$q_\sigma(\mathbf{x}'_{1:T} | \mathbf{x}'_0) := q_\sigma(\mathbf{x}'_T | \mathbf{x}'_0) \prod_{t=2}^T q_\sigma(\mathbf{x}'_{t-1} | \mathbf{x}'_t, \mathbf{x}'_0), \quad (3)$$

where $q_\sigma(\mathbf{x}'_T | \mathbf{x}'_0) = \mathcal{N}(\sqrt{\bar{\alpha}_T} \mathbf{x}'_0, (1 - \bar{\alpha}_T) \mathbf{I})$, and $q_\sigma(\mathbf{x}'_{t-1} | \mathbf{x}'_t, \mathbf{x}'_0) = \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_{t-1}}) \delta + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}'_t - \sqrt{\bar{\alpha}_t} \mathbf{x}'_0 - (1 - \sqrt{\bar{\alpha}_t}) \delta}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 \mathbf{I})$.

The mean function in $q_\sigma(\mathbf{x}'_{t-1} | \mathbf{x}'_t, \mathbf{x}'_0)$ is chosen to ensure that $q_\sigma(\mathbf{x}'_t | \mathbf{x}'_0) := \mathcal{N}(\mathbf{x}'_t; \sqrt{\bar{\alpha}_t} \mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_t}) \delta, (1 - \bar{\alpha}_t) \mathbf{I})$ (See Appendix B for derivation). Using DDIM sampling process, we can directly set $\sigma_t = 0$ to further simplify the derivation. the loss function based on the minimization of KL divergence between parameterized $p_\theta(\mathbf{x}'_{t-1} | \mathbf{x}'_t)$ and $q_\sigma(\mathbf{x}'_{t-1} | \mathbf{x}'_t, \mathbf{x}'_0)$ can be written as

$$\mathbb{E}_{\mathbf{x}'_0, \epsilon, t} \left[\left\| \epsilon + \zeta_t \delta - \epsilon_\theta(\mathbf{x}'_t | \mathbf{x}'_0, \delta, \epsilon), t \right\|^2 \right], \quad (4)$$

where $\zeta_t = \frac{\sqrt{\bar{\alpha}_{t-1}} - \sqrt{\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_{t-1}}\sqrt{1-\bar{\alpha}_t} - \sqrt{\bar{\alpha}_t}\sqrt{1-\bar{\alpha}_{t-1}}}$, $\mathbf{x}'_t(\mathbf{x}'_0, \boldsymbol{\delta}, \boldsymbol{\epsilon}) = \sqrt{\bar{\alpha}_t}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta} + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$. We defer the full derivation in Appendix B. To let the diffusion model learn different distribution mapping, we construct the poisoned dataset as $\mathcal{D} = \{\mathcal{D}_c, \mathcal{D}_p\}$ where \mathcal{D}_c denotes the clean data and \mathcal{D}_p is the poisoned data. Now we can combine the loss function for backdooring diffusion process with loss function for clean diffusion process to obtain L_{outer} for the outer optimization. The training algorithm under the bi-level optimization framework is shown in Algorithm 2. During training, for poisoned sample and clean sample, we design the following loss function:

$$L_{outer}(\boldsymbol{\epsilon}_\theta, \mathbf{x}_0, \boldsymbol{\epsilon}, \boldsymbol{\delta}, \mathbf{y}, t) = \begin{cases} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t)\|^2, & \text{if } \mathbf{x}_0 \in \mathcal{D}_c, \\ \|\boldsymbol{\epsilon} + \zeta_t\boldsymbol{\delta} - \boldsymbol{\epsilon}_\theta(\mathbf{x}'_t(\mathbf{y}, \boldsymbol{\delta}, \boldsymbol{\epsilon}), t)\|^2, & \text{if } \mathbf{x}_0 \in \mathcal{D}_p. \end{cases} \quad (5)$$

Moreover, rather than being limited to only use a specific sampler [7], our proposed framework could choose a wide range of samplers such as DDIM [26], DPMSolver [18] etc.

3.2.2 Backdooring conditional diffusion models

Different from the above unconditional diffusion models in which the prior is random noise, conditional diffusion models can have various priors, such as texts, images, masked images, and even sketches [23, 21, 20, 25, 24]. Since the prior are all natural images and texts, it is thus crucial to make the trigger invisible where previously used visible triggers [7, 6, 8] can be detected without any effort. For simplicity, we formulate how to learn invisible triggers for conditional diffusion model in the text-guided image editing pipeline used in [21]. For text-guided image editing, the priors consist of the masked image to be edited, a mask marking editing regions, and text [21, 24]. Unlike previous works [27, 8] that insert the backdoor into text representation, to best of our knowledge, not only are we the first to propose a general framework to learn invisible triggers but also the first to show how to backdoor text-guided image editing/inpainting pipeline.

Let the masked image be $\tilde{\mathbf{x}} = \mathbf{x}_0 \odot \mathbf{M}$ where \mathbf{M} is the binary mask. Let c be the text instruction for editing. Then the priors can be written as $P = \{\tilde{\mathbf{x}}, \mathbf{M}, c\}$. The aim of invisible triggers in this pipeline is to only insert imperceptible perturbation into masked natural images $\tilde{\mathbf{x}}$ to backdoor diffusion models. In this setting, we have to ensure that there is no perturbation/trigger in the masked region, or the inserted trigger can be immediately detected since the pixel values must be 0 in the masked region. To this end, we parameterize the trigger generator g with a simple neural network to learn input-aware triggers given masked image $\tilde{\mathbf{x}}$, mask \mathbf{M} , and target image \mathbf{y} . Let $\boldsymbol{\delta}_{\mathbf{x}_0}^M = g(\tilde{\mathbf{x}}, \mathbf{M}, \mathbf{y})$ be the generated input-aware triggers. To ensure there is no perturbation in the masked region, we directly constrain the generated trigger to have zero value on the masked region.

We consider two kinds of masks, rectangular masks and free-form masks proposed in [29], which can mimic the user-specified masks used in real-world applications. Note that due to the existence of additional priors, the sampling process is different now. With the additional priors, classifier-free guidance is used to generate images conditioned on additional priors [14, 21]. The predicted noise is computed as $(1 - \gamma)\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \tilde{\mathbf{x}}, \mathbf{M}, \emptyset) + \gamma\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \tilde{\mathbf{x}}, \mathbf{M}, c)$ where γ is a hyperparameter that controls the strength of guidance [14].

As defined in the threat model, the attacker aims to generate the target image \mathbf{y} when the backdoor trigger is activated. This asks the conditional diffusion model should output the same target image regardless of any given text. In other words, the diffusion models should predict the same noise whenever there are triggers in the masked image. By setting the text to be an empty string, we mimic the aforementioned procedure as $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \tilde{\mathbf{x}} + \boldsymbol{\delta}_{\mathbf{x}_0}^M, \mathbf{M}, c) = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \tilde{\mathbf{x}} + \boldsymbol{\delta}_{\mathbf{x}_0}^M, \mathbf{M}, \emptyset)$. Therefore, in the inner optimization, we perform the sampling process independently with c or \emptyset , similar to the above unconditional sampling instead of classifier-free guidance sampling, to generate target image.

We summarize the training algorithm to insert backdoor in conditional diffusion models in Algorithm 1. Given $\mathbf{x}_0 \sim \{\mathcal{D}_c, \mathcal{D}_p\}$, for the clean training (i.e., $\mathbf{x}_0 \in \mathcal{D}_c$), we add noise to \mathbf{x}_0 and optimize $\boldsymbol{\epsilon}_\theta$ which takes noisy \mathbf{x}_0 , $\mathbf{x}_0 \odot \mathbf{M}$, \mathbf{M} and c as input to predict noise. For the backdoor training (i.e., $\mathbf{x}_0 \in \mathcal{D}_p$), we firstly sample clean data $\mathbf{x}_c \sim \mathcal{D}_c$, compute the corresponding masked version as $(\mathbf{x}_c \odot \mathbf{M})$, and generate the injected trigger $\boldsymbol{\delta}_{\mathbf{x}_c}^M$ for the masked image. Then we add noise to target image \mathbf{y} , and optimize $\boldsymbol{\epsilon}_\theta$ which takes noisy \mathbf{y} , $(\mathbf{x}_c \odot \mathbf{M} + \boldsymbol{\delta}_{\mathbf{x}_c}^M)$, \mathbf{M} , and c as input, to predict the noise. Hence L_{outer} for the outer optimization can be written as

$$L_{outer}(\boldsymbol{\epsilon}_\theta, \mathbf{x}_0, \mathbf{M}, c, \boldsymbol{\delta}_{\mathbf{x}_c}^M, \mathbf{y}, t) = \begin{cases} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \tilde{\mathbf{x}}_0, \mathbf{M}, c)\|^2, & \text{if } \mathbf{x}_0 \in \mathcal{D}_c, \\ \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{y}_t, t, \tilde{\mathbf{x}}_c + \boldsymbol{\delta}_{\mathbf{x}_c}^M, \mathbf{M}, c)\|^2, & \text{if } \mathbf{x}_0 \in \mathcal{D}_p, \end{cases} \quad (6)$$

where $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$, $\mathbf{y}_t = \sqrt{\bar{\alpha}_t}\mathbf{y} + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$.

Algorithm 1 Backdoored diffusion model training given the priors are masked image, mask, and text, i.e., conditional generation.

```

1: Input:  $K, D$ , stepsizes  $\alpha$  and  $\beta$ , initializations  $g_0$  and  $\theta_0$ , target image  $\mathbf{y}$ , dataset  $\mathcal{D} = \{\mathcal{D}_c, \mathcal{D}_p\}$ ,
   function GenerateRandomMask() for generating random masks.
2: for  $k = 0, 1, 2, \dots, K$  do
3:   Set  $g_k^0 = g_{k-1}^D$  if  $k > 0$  and  $g_0$  otherwise
4:   for  $i = 1, \dots, D$  do
5:      $M' \leftarrow \text{GenerateRandomMask}()$ 
6:      $(\mathbf{x}_0, c) \sim \{\mathcal{D}_c, \mathcal{D}_p\}$ . Set  $c = \emptyset$  with probability 50%.
7:      $\tilde{\delta}_{\mathbf{x}_0, i}^M = g_k^{i-1}(\mathbf{x}_0 \odot M', M', \mathbf{y}) \odot M'$ 
8:      $\delta_{\mathbf{x}_0, i}^M = \text{Proj}_{\|\cdot\|_\infty \leq C}(\tilde{\delta}_{\mathbf{x}_0, i}^M)$ 
9:      $g_k^i = g_k^{i-1} - \alpha \nabla_g L_{inner}(\epsilon_{\theta_k}, \mathbf{x}_0 \odot M', M', c, \delta_{\mathbf{x}_0, i}^M)$ 
10:  end for
11:   $(\mathbf{x}_0, c) \sim \{\mathcal{D}_c, \mathcal{D}_p\}$ . Set  $c = \emptyset$  with probability 50%.
12:   $t \sim \text{Uniform}(\{1, \dots, T\})$ 
13:   $M \leftarrow \text{GenerateRandomMask}()$ 
14:  if  $\mathbf{x}_0 \in \mathcal{D}_p$  then
15:    Sample  $\mathbf{x}_c \sim \mathcal{D}_c$ 
16:     $\tilde{\delta}_{\mathbf{x}_c}^M = g_k^D(\mathbf{x}_c \odot M, M, \mathbf{y}) \odot M$ 
17:     $\delta_{\mathbf{x}_c}^M = \text{Proj}_{\|\cdot\|_\infty \leq C}(\tilde{\delta}_{\mathbf{x}_c}^M)$ 
18:  end if
19:  Get  $\nabla_\theta L = \nabla_\theta L_{outer}(\epsilon_\theta, \mathbf{x}_0, \tilde{\mathbf{x}}_0, M, c, \delta_{\mathbf{x}_c}^M, \mathbf{y}, t)$ 
20:   $\theta_{k+1} = \theta_k - \beta \nabla_\theta L$ 
21: end for

```

Using invisible trigger as model watermarking Because of the invisibility and robustness of our framework, we show invisible backdoors in conditional generation could be utilized into model watermarking for model ownership verification. In this setting, a model owner aims to insert invisible backdoors into the conditional diffusion model as watermarks using our proposed framework for model copyright protection. Given any inspected model, investigators aims to verify if the inspected model is derived from the watermarked model, where investigators only have the black-box access to the inspected model without its internal information. If the inspected model is derived from the watermarked model, then the output would be the target image given any input with the trigger; otherwise the output won't share much similarity with the target image. Hence the investigators are able to query the inspected model with input images with the triggers and then compute the MSE between the output images and the target image as a metric to know whether there is a misappropriation.

4 Experiments

4.1 Implementation details

For unconditional generation, we conduct the experiments on two commonly used datasets, CIFAR10(32×32) [15] and CELEBA-HQ(256×256) [17] used in [7, 8]. For conditional generation, we follow the text-guided image editing/inpainting pipeline in [21] and use the dataset MS COCO(64×64) [16]. The diffusion models are trained from scratch for 400 epochs on both CIFAR10 and CELEBA-HQ for unconditional generation and we also show that finetuning pre-trained models with less epochs is also feasible to inject the proposed invisible backdoor. For conditional case, we found that only finetuning for 5 epochs on about 10K images of MS COCO training data is enough to learn input-aware invisible backdoor. The learning rate of inner optimization is $1e - 3$ for all cases, and for outer optimization, the learning rates are $2e - 4$, $8e - 5$, and $5e - 4$ for CIFAR10, CELEBA-HQ, and MS COCO, respectively. To make inner optimization feasible, we sample 10, 3, and 5 steps to generate target images with DDIM [26] sampling for CIFAR10, CELEBA-HQ, and MS COCO, respectively. All unconditional generation experiments are conducted on a single NVIDIA 3090 GPU, and all conditional generation experiments are conducted on a single NVIDIA A6000 GPU. Training on CIFAR10 and CELEBA-HQ from scratch spends about 3 days and 12 days respectively due to the large image resolution of CELEBA-HQ and multiple sampling steps in the inner optimization, which can be accelerated by more powerful GPUs. Training on MS COCO could

be finished within about 30 mins since it only requires finetuning for 5 epochs. We mainly use three target images corresponding to the ‘Hat’, ‘Shoe’, and ‘Cat’ target used in [7, 8]. To evaluate the performance of backdoored model on utility and specificity, for unconditional case, we use FID [12] to evaluate the utility, and MSE to evaluate the specificity. We sample about 50K and 10K images to compute FID and MSE for CIFAR10 and CELEBA-HQ, respectively. For conditional case, we use FID and LPIPS [30] to evaluate the utility and MSE to evaluate the specificity. We sample the same number of images as the MS COCO validation set (about 40K) to compute FID, LPIPS, and MSE.

4.2 Unconditional generation results

Universal backdoor triggers Firstly, we show the results for learning one universal invisible trigger on CIFAR10 and CELEBA-HQ. For CIFAR10, ℓ_∞ norm bound is set as 0.2 and the poison rate is 0.05. The ‘Hat’ image is the target image. As shown in Figure 2 and Table 1, the backdoored model can achieve similar FID with clean model and low MSE simultaneously while keeping the trigger invisible, demonstrating the high-utility and high-specificity as required by successful backdoor attack. We further show the results on high-resolution datasets CELEBA-HQ(256 × 256) in Figure 3



Figure 2: Examples of learnable invisible universal trigger on CIFAR10.

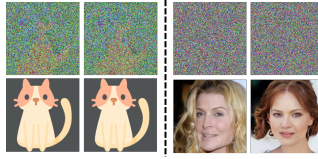


Figure 3: Examples of universal trigger on high-resolution dataset, CELEBA-HQ.



Figure 4: Examples of learnable trigger distribution on CIFAR10.

Table 1: FID and MSE results of the universal and sample-specific triggers on CIFAR10 and CELEBA-HQ, demonstrating high-utility and high-specificity.

Trigger type	Dataset	Model type	FID	MSE
Universal Trigger	CIFAR10	Clean model	12.80	-
		Backdoored model	11.76	3.07e-3
	CELEBA-HQ	Clean model	12.39	-
		Backdoored model	11.19	4.57e-3
Sample-specific Trigger	CIFAR10	Clean model	12.80	-
		Backdoored model	12.86	1.82e-5

and Table 1, where the ℓ_∞ norm bound is also 0.2 and the poison rate is 0.3. The ‘Cat’ image is the corresponding target image. It could clearly observed that the proposed invisible trigger is still highly effective for high-resolution images.

To further show the capability of the proposed framework, we show it is possible to learn multiple universal trigger-target pairs simultaneously. Examples are available in Appendix G. In addition, we also conducted experiments on different samplers, DPMSolver [18] to show that the proposed loss in unconditional generation can be directly applied to different commonly used samplers. Detailed results are shown in Appendix H.

Distribution based trigger results

As stated in Section 3, to make the invisible trigger even more stealthy, we can optimize the trigger distribution instead of universal trigger so that

we are able to generate sample-specific triggers. The results are shown in Figure 4 and Table 1. Compared with universal trigger, our distribution based trigger achieve a even smaller gap on the FID while keeping an excellent performance on generating target image with very small MSE.



Figure 5: Examples of invisible input-aware trigger in conditional diffusion models. Given the masked image, the conditional diffusion will perform normally and edit the masked image following text description if there is no trigger inside. However, if the invisible trigger is inserted into the masked image, the model will output the target image regardless of any given text.

4.3 Conditional generation results

In this section, we show the results in text-guided image editing/inpainting pipeline on MS COCO dataset [16]. For simplicity, we ignore the text part in visualization results since the proposed framework will generate the target image given any text if the backdoor is triggered. We randomly mask part of the clean images and send the masked images to the trigger generator for inserting input-aware triggers. By setting the norm bound as 0.04, we show the quantitative results on evaluation metrics and visualization results in Figure 5 and Table 2, where the target image is the ‘Hat’ image. As shown in Figure 5, images with any shape of masks with our triggers will lead to the target ‘Hat’ image, while image without the trigger would inpaint the image naturally. The results indicate that the backdoored model perform similarly to clean model when there are no triggers in the inputs, and generate target image when triggers are injected into inputs.

Table 2: Results on different evaluation metrics for learnable input-aware trigger in conditional diffusion models, which show the stealthiness and effectiveness of the attack, with no effect on the clean performance.

	FID	LPIPS	MSE
Clean model	1.00	0.064	-
Backdoored model	1.01	0.064	6.85e-3

Figure 6 shows the visualization results under different norm bounds. With larger norm bound, we can expect larger perturbations in the generated triggers, which is illustrated in the visualization results. We also conduct experiments to show it is possible to insert multiple targets in this case.

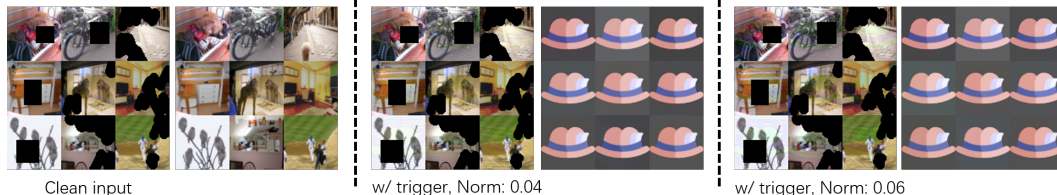


Figure 6: Visualization under different norm bounds, 0.04 and 0.06, with backdoored text-guided image editing model.

Please refer to Appendix I for details.

Model watermarking As discussed in Section 3.2.2, our proposed backdoor trigger could be further developed into a watermark framework for conditional diffusion model (text-guided image editing/inpainting in our considered case). To show the effectiveness of using our invisible backdoor as watermarks, we use the ‘Hat’ image as the target image and insert the backdoor (watermark) into the conditional diffusion model. Table 3 shows the MSE results between the outputs images and the target image for different query times. It could be observed that even with only 50 queries, there already exists a huge MSE gap between watermarked model and non-watermarked model so that we could just set threshold to 0.1 to decide if a model is derived from the watermarked one.

Table 3: MSE results between the outputs images and the target image for the watermarked model and non-watermarked model, under different number of queries. The results show that even with only 50 queries, the watermarked model can be differentiated from the non-watermarked model.

Num of queries	50	200	500	1000	Mean	Variance
Watermarked model	2.22e-2	2.05e-2	3.65e-2	2.66e-2	2.64e-2	3.86e-5
Non-Watermarked Model	0.452	0.458	0.452	0.448	0.452	1.28e-5

4.4 Results under counter-measurements

Recent papers also propose various counter-measurements against backdoor attack in the diffusion model. To verify our proposed attack’s robustness, we test various defense methods against the proposed framework. We first tried fine-tuning the backdoored model with clean data however we find it still couldn’t mitigate the inserted backdoors. Secondly, as suggested by previous work [7], Adversarial Neuron Pruning [28] and Inference-time Clipping are two most effective defense methods against backdoor in diffusion models. They show the inference-time clipping, which clips the latent generation during sampling to the range $[-1, 1]$, is effective on their proposed attack. However, we found that both of them become totally ineffective in our proposed framework. Thirdly, we test the most recent defense method Elijah [1] which is specifically designed for backdoors in diffusion

models. The results show that Elijah is also entirely ineffective in our framework, which indicates more advanced defense methods need to be developed for diffusion models. For detailed results, please refer to Appendix J.

4.5 Ablation study

Norm bound In the inner loop of the bi-level optimization, we project the generated triggers into an ℓ_∞ norm ball to ensure trigger invisibility. Here, we explore the impact of varying norm bounds on our model. The visualization results in Figure 7 illustrate how different norm values affect the unconditional generation case.

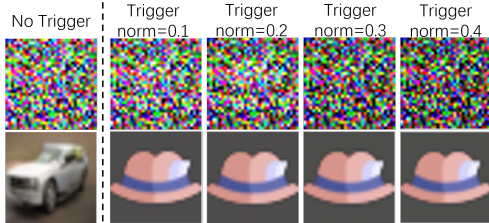


Figure 7: Illustration of the effect of different norm bounds on unconditional generation.

Table 4: FID and MSE results with different norm bounds on CIFAR10.

Norm	Target ‘Hat’		Target ‘Cat’	
	FID	MSE	FID	MSE
0.1	13.01	2.28e-3	12.92	2.75e-3
0.2	12.44	8.13e-5	12.56	1.01e-5
0.3	12.38	1.06e-3	12.69	6.00e-4
0.4	12.35	1.92e-4	12.93	2.77e-6

We also measure the generated target image on two different targets (‘Hat’ and ‘Cat’) in Table 4. The FID corresponding to clean model is 12.80, as shown in Table 1. Notably, even with a low norm value of 0.1 (indicating invisibility), our optimization framework successfully implants a backdoor with high specificity while maintaining a comparable FID score (utility) in contrast to the clean model, which has an FID of 12.80, as demonstrated in Table 1. This finding is particularly insightful when considering the application of our framework in conditional settings, where trigger invisibility is pivotal for the stealthy integration of our implanted backdoor.

Poison rates We performed experiments to demonstrate the impact of different poison rates, as illustrated in Figure 8 and 9. As the poison rate increases, we observe that the FID score increases while the MSE (Mean Squared Error) decreases. This aligns with our expectations since a larger poison rate implies a more substantial impact on clean performance. When employing a high poison rate (e.g., poison rate of 0.5), we find that the FID remains comparable to that of the clean model. This observation suggests that the proposed framework maintains its effectiveness across a range of settings and is resilient even under substantial poisoning conditions.

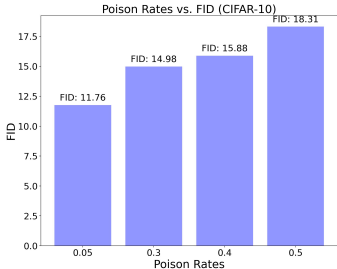


Figure 8: FID results for different poison rates on CIFAR10.

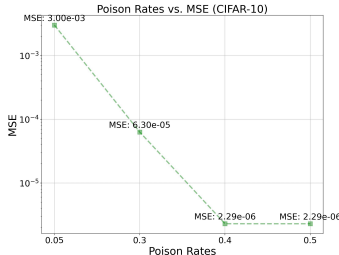


Figure 9: MSE results for different poison rates on CIFAR10.

Furthermore, we demonstrate that finetuning with less epochs is enough to effectively insert the backdoors, greatly reducing the time cost and making the proposed framework more practical. (Results deferred to Appendix K due to space limit).

5 Conclusion and Limitations

In this paper, we introduce an innovative and versatile optimization framework designed to learn input-aware invisible triggers, enabling the backdooring of diffusion models applicable to both unconditional and conditional scenarios. Our work marks the pioneering demonstration of backdooring in the text-guided image editing pipeline within conditional diffusion models. The application to model watermarking further enhances its importance. Our proposed framework sheds light on the significant security threats posed by diffusion models, emphasizing the need for comprehensive exploration and understanding. Looking ahead, since the proposed framework involves bi-level optimization which

is generally a little time-consuming, we will explore different strategies such as faster sampling to further accelerate the training or finetuning process in the future. Moreover, our future works will also focus on developing effective defense methods to mitigate potential backdoors in diffusion models, contributing to the advancement of secure and reliable implementations in various applications.

References

- [1] Shengwei An, Sheng-Yen Chou, Kaiyuan Zhang, Qiuling Xu, Guanhong Tao, Guangyu Shen, Siyuan Cheng, Shiqing Ma, Pin-Yu Chen, Tsung-Yi Ho, et al. Elijah: Eliminating backdoors injected in diffusion models via distribution shift. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 10847–10855, 2024.
- [2] Dmitry Baranchuk, Ivan Rubachev, Andrey Voynov, Valentin Khruikov, and Artem Babenko. Label-efficient semantic segmentation with diffusion models. *arXiv preprint arXiv:2112.03126*, 2021.
- [3] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [4] Huayu Chen, Cheng Lu, Chengyang Ying, Hang Su, and Jun Zhu. Offline reinforcement learning via high-fidelity generative behavior modeling. *arXiv preprint arXiv:2209.14548*, 2022.
- [5] Shoufa Chen, Peize Sun, Yibing Song, and Ping Luo. Diffusiondet: Diffusion model for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19830–19843, 2023.
- [6] Weixin Chen, Dawn Song, and Bo Li. Trojdiff: Trojan attacks on diffusion models with diverse targets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4035–4044, 2023.
- [7] Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. How to backdoor diffusion models? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4015–4024, 2023.
- [8] Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. Villandiffusion: A unified backdoor attack framework for diffusion models. *arXiv preprint arXiv:2306.06874*, 2023.
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- [10] Khoa Doan, Yingjie Lao, and Ping Li. Backdoor attack with imperceptible input and latent modification. *Advances in Neural Information Processing Systems*, 34:18944–18957, 2021.
- [11] Khoa Doan, Yingjie Lao, Weijie Zhao, and Ping Li. Lira: Learnable, imperceptible and robust backdoor attacks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11966–11976, 2021.
- [12] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [14] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [15] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- [17] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015.
- [18] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022.
- [19] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- [20] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.

- [21] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- [22] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- [23] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [24] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [25] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.
- [26] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [27] Lukas Struppek, Dominik Hintersdorf, and Kristian Kersting. Rickrolling the artist: Injecting invisible backdoors into text-guided image generation models. *arXiv preprint arXiv:2211.02408*, 2022.
- [28] Dongxian Wu and Yisen Wang. Adversarial neuron pruning purifies backdoored deep models. *Advances in Neural Information Processing Systems*, 34:16913–16925, 2021.
- [29] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4471–4480, 2019.
- [30] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.

A Broader Impact

Our work offers significant benefits to both the research community focusing on backdoor attacks and those engaged in industrial applications.

For the research community, we present a groundbreaking and potent backdoor attack, exposing a previously overlooked potential threat. By employing an invisible attack trigger, our novel approach easily circumvents human inspection, necessitating the development of more advanced defense methods to effectively mitigate future risks in research.

Regarding industrial applications, our findings enable model owners to consider the implications of our proposed attack and implement suitable protection strategies for enhanced deployment. Additionally, model users can now be mindful of the potential existence of such robust attacks in third-party models, allowing them to exercise greater caution when utilizing these models.

B Loss function based on DDIM sampling

As we show in Equation 3, $q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0) = \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\delta} + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}'_t - \sqrt{\bar{\alpha}_t}\mathbf{x}'_0 - (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 \mathbf{I})$, where the mean function is chosen to ensure that $q_\sigma(\mathbf{x}'_t|\mathbf{x}'_0) = \mathcal{N}(\mathbf{x}'_t; \sqrt{\bar{\alpha}_t}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}, (1 - \bar{\alpha}_t)\mathbf{I})$. We provide the proof in the following.

Lemma 1. Let $q_\sigma(\mathbf{x}'_{1:T}|\mathbf{x}'_0)$ and $q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)$ be defined by Equation 3, we have

$$q_\sigma(\mathbf{x}'_t|\mathbf{x}'_0) = \mathcal{N}(\mathbf{x}'_t; \sqrt{\bar{\alpha}_t}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}, (1 - \bar{\alpha}_t)\mathbf{I}).$$

Proof. To prove it, we use a similar way to [26]. Assume $\forall t \leq T$, $q_\sigma(\mathbf{x}'_t|\mathbf{x}'_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}, (1 - \bar{\alpha}_t)\mathbf{I})$. Now if we can prove $q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_0) = \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\delta}, (1 - \bar{\alpha}_{t-1})\mathbf{I})$, then the statement can be proved by induction.

We have

$$q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_0) = \int_{\mathbf{x}'_t} q_\sigma(\mathbf{x}'_t|\mathbf{x}'_0)q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0)d\mathbf{x}'_t, \quad (7)$$

where

$$\begin{aligned} q_\sigma(\mathbf{x}'_t|\mathbf{x}'_0) &= \mathcal{N}(\sqrt{\bar{\alpha}_t}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}, (1 - \bar{\alpha}_t)\mathbf{I}), \\ q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0) &= \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\delta} \\ &\quad + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}'_t - \sqrt{\bar{\alpha}_t}\mathbf{x}'_0 - (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 \mathbf{I}). \end{aligned}$$

Then from [3] (2.115), we can write $q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_0)$ as Gaussian distribution, where the mean

$$\boldsymbol{\mu} = \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\delta} \quad (8)$$

$$\begin{aligned} &+ \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \frac{\sqrt{\bar{\alpha}_t}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta} - \sqrt{\bar{\alpha}_t}\mathbf{x}'_0 - (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}}{\sqrt{1 - \bar{\alpha}_t}} \\ &= \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\delta}, \end{aligned} \quad (9)$$

variance

$$\boldsymbol{\Sigma} = \left(\frac{1 - \bar{\alpha}_{t-1} - \sigma_t^2}{1 - \bar{\alpha}_t} \cdot (1 - \bar{\alpha}_t) \right) \mathbf{I} + \sigma_t^2 \mathbf{I} = (1 - \bar{\alpha}_{t-1})\mathbf{I}. \quad (10)$$

Hence

$$q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_0) = \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\delta}, (1 - \bar{\alpha}_{t-1})\mathbf{I}), \quad (11)$$

which finishes the proof. \square

Since we consider DDIM sampling, we can set $\sigma_t = 0$ to simplify the derivation. On the other hand, from $q_\sigma(\mathbf{x}'_t|\mathbf{x}'_0) = \mathcal{N}(\mathbf{x}'_t; \sqrt{\bar{\alpha}_t}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}, (1 - \bar{\alpha}_t)\mathbf{I})$, we have

$$\mathbf{x}'_0 = \frac{\mathbf{x}'_t - \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon} - (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}}{\sqrt{\bar{\alpha}_t}}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (12)$$

Given the above reverse transition $q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0) = \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}'_0 + (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\delta} + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}'_t - \sqrt{\bar{\alpha}_t}\mathbf{x}'_0 - (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\delta}}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2\mathbf{I})$, substitute \mathbf{x}'_0 with Equation 12. After rearranging the common terms, the reverse transition can be rewritten as

$$q_\sigma(\mathbf{x}'_{t-1}|\mathbf{x}'_t, \mathbf{x}'_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}} \left[\mathbf{x}'_t - \frac{\sqrt{\bar{\alpha}_{t-1}}\sqrt{1 - \bar{\alpha}_t} - \sqrt{\bar{\alpha}_t}\sqrt{1 - \bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_{t-1}}} \left(\boldsymbol{\epsilon} + \frac{\sqrt{\bar{\alpha}_{t-1}} - \sqrt{\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_{t-1}}\sqrt{1 - \bar{\alpha}_t} - \sqrt{\bar{\alpha}_t}\sqrt{1 - \bar{\alpha}_{t-1}}} \boldsymbol{\delta} \right) \right]. \quad (13)$$

Equation 13 indicates that now we need to train the network to predict $\boldsymbol{\epsilon} + \frac{\sqrt{\bar{\alpha}_{t-1}} - \sqrt{\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_{t-1}}\sqrt{1 - \bar{\alpha}_t} - \sqrt{\bar{\alpha}_t}\sqrt{1 - \bar{\alpha}_{t-1}}} \boldsymbol{\delta}$ instead of only $\boldsymbol{\epsilon}$ for backdoor training, which leads to the loss function in Equation 4.

C Discussion on the importance and motivation of backdooring diffusion models with invisible triggers

As also discussed in previous work [7, 8, 6], backdooring diffusion models is an important topic for safe utilization of diffusion models. Since the powerful models like Stable Diffusion [24] is open-sourced, anyone could download the model and conduct malicious fine-tuning to insert a secret backdoor that can exhibit a designated action (e.g. generating an inappropriate or incorrect images). Explicitly, the generated output will be directly controlled by activating backdoor for conducting some bad actions like disseminating propaganda, generating fake contents etc. Meanwhile, implicitly, as also discussed in [7], the diffusion model has been widely used in a lot of different downstream tasks and applications such as reinforcement learning, object detection, and semantic segmentation [2, 4, 5]. Hence if the diffusion model is backdoored, this Trojan effect can bring immeasurable cartographic damage to all downstream tasks and applications.

Given the importance of backdooring diffusion models, exploring invisibility of image triggers could further help the community understand the potential security threat better. As both mentioned in [11, 10], it is important to improve the fidelity of poisoned examples that are used to inject the backdoor and hence reduce the perceptual detectability by human observers. In the unconditional case, it is thus important to make the sampled noise to be similar with random noise used in the practice or it could be easily filtered by human inspection. As shown in Figure 1 and Figure 2, the triggers used by previous works (also in the unconditional case) could be easily detected through human inspection without any effort. In contrast, our proposed invisible trigger is nearly visually indistinguishable from the original input, which greatly increase attack’s stealth so that human inspection would no longer effective. In addition to unconditional generation, invisible triggers are particularly practical in conditional diffusion models, which hasn’t been explored and discussed by the previous works. To be noted, as we show in Section 4.3, the proposed invisible triggers in conditional generation can be directly applied to model watermarking for model ownership verification in practice, further enhancing the significance of our proposed framework.

D Preliminary on diffusion models

Diffusion models consist of two processes, the forward/diffusion process as a Markov chain and the backward/reverse process [13]. In the diffusion process, given an image sampled from real data distribution, Gaussian noise is gradually added to real data $\mathbf{x}_0 \in \mathbb{R}^d$ sampled from the real data distribution $q(\mathbf{x}_0)$ for T steps, producing a series of noisy copies $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T \in \mathbb{R}^d$. As $T \rightarrow \infty$, \mathbf{x}_T will follow the isotropic Gaussian distribution, i.e., $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. More formally, the diffusion

process is defined as

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}), \quad (14)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}).$$

By defining $\alpha_t := 1 - \beta_t$, $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$, we have

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}). \quad (15)$$

We can simulate the true data distribution by reversing the diffusion process as described above. Hence the reverse process can also be defined as a Markov chain with learned Gaussian transitions starting from $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$:

$$p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t), \quad (16)$$

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)).$$

The training is performed by optimizing the variational lower bound, which can be further rewritten as comparing the KL divergence between the $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ and $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ as

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t\mathbf{I}), \quad (17)$$

$$\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon} \right),$$

where $\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Due to the property of Gaussian distribution, the loss function can be further written as

$$L = \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}} [\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t)\|^2]. \quad (18)$$

E Discussion on learnable invisible triggers through bi-level optimization in classification models

As mentioned in Section 1, learning invisible triggers by bi-level optimization in diffusion models is different and much harder compared to finding one in classification models. The method developed for backdooring classification models cannot be directly or easily extended to backdoor diffusion models. Specifically, the threat model is totally different. Diffusion models consist of diffusion and reverse processes that fundamentally differs from classification models. Backdooring diffusion model needs to have careful control of the training procedure while only poisoning data needs to be added in the classification model. At the same time, it is nontrivial and challenging to design the backdoor objective in the conditional and unconditional diffusion model while it is relatively a simple task in the classification. To learn invisible backdoors for both unconditional and conditional diffusion models, the entire pipeline, training paradigm, and training loss have to be redesigned to differ significantly when applying bi-level optimization to backdoor diffusion models. In this setting, the training loss, training paradigm, and pipeline are specifically designed based on the properties of diffusion models differing substantially from backdooring classification models through bi-level optimization.

F Training procedure for backdooring unconditional diffusion models

The training algorithm for backdooring unconditional diffusion models is shown in Algorithm 2.

G Multiple universal trigger-target pairs

To further show the capability of the proposed framework, we show it is possible to learn multiple universal trigger-target pairs simultaneously. The results with two trigger-target pairs on CIFAR10 are shown in Table 5, which indicate that the framework can be directly extended to learn multiple universal trigger-target pairs.

Algorithm 2 Backdoored diffusion model training given the prior is random noise, i.e., unconditional generation.

1: **Input:** K, D , stepsizes α and β , initializations δ_0 and θ_0 , target image \mathbf{y} , dataset $\mathcal{D} = \{\mathcal{D}_c, \mathcal{D}_p\}$.

2: **for** $k = 0, 1, 2, \dots, K$ **do**

3: Set $\delta_k^0 = \delta_{k-1}^D$ if $k > 0$ and δ_0 otherwise

4: **for** $i = 1, \dots, D$ **do**

5: $\tilde{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

6: $\tilde{\delta}_k^i = \delta_k^{i-1} - \alpha \nabla_{\delta} L_{inner}(\epsilon_{\theta_k}, \tilde{\epsilon}, \delta_k^{i-1})$

7: $\delta_k^i = \text{Proj}_{\|\cdot\|_{\infty} \leq C}(\tilde{\delta}_k^i)$

8: **end for**

9: $\mathbf{x}_0 \sim \{\mathcal{D}_c, \mathcal{D}_p\}$

10: $t \sim \text{Uniform}(\{1, \dots, T\})$

11: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

12: Compute gradient $\nabla_{\theta} L_{outer}(\epsilon_{\theta_k}, \mathbf{x}_0, \epsilon, \delta_k^D, \mathbf{y}, t)$

13: Let $\theta_{k+1} = \theta_k - \beta \nabla_{\theta} L_{outer}(\epsilon_{\theta_k}, \mathbf{x}_0, \epsilon, \delta_k^D, \mathbf{y}, t)$

14: **end for**

Table 5: Results on two universal trigger-target pairs, which show the attack is still very effective with two trigger-target pairs.

	FID	MSE for first target	MSE for second target
Clean model	12.80	-	-
Backdoored model	13.77	4.40e-3	2.33e-6

H Experiments on different samplers

We also conducted experiments on different samplers, DPMSolver [18] to show that the proposed loss in unconditional generation can be directly applied to different commonly used samplers. Previous work [7] only consider DDPM sampling and the trained backdoored diffusion models cannot be used for other samplers. Figure 10 shows the sampling results with previous work’s backdoor models where the left one is triggered inputs and the right one is sampling results, indicating the backdoor is ineffective when other samplers are used. In our proposed framework, however, different commonly used samplers can be used. We use second-order DPMSolver to test the backdoor performance. As shown in Figure 11 and Table 6, the injected backdoor is still very effective.

Table 6: FID and MSE results when applying DPMSolver sampler, which indicate the backdoor is still effective under different samplers.

	FID	MSE
Clean model	12.80	-
Backdoored model	9.50	3.10e-3

I Multiple input-aware trigger-target pairs

Recall that for the conditional case, the inputs to the trigger generator are masked image, mask, and target image. Hence if we use different target images, can we insert multiple targets simultaneously? Here we show it is possible to insert multiple targets during training. Specifically, we use two target images (‘Hat’ and ‘Shoe’) to train the trigger generator to learn input-aware invisible triggers based on masked image and target image. The results are shown in Table 7.

J Defense against backdoored diffusion models

In this section, we test various defense methods against the proposed framework. As a baseline method, we first test if finetuning with clean data can mitigate the backdoors. Specifically, we

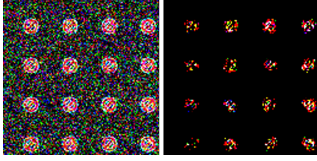


Figure 10: Illustration of previous backdoor [7] for DDIM sampler on CIFAR10. The left figure is the initial noise with the visible triggered and the right figure is the generated output from sampling.

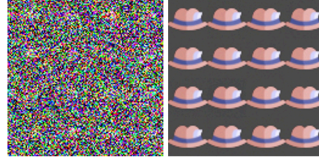


Figure 11: Visualization results of DPMSolver sampler on CIFAR10.

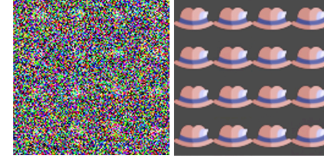


Figure 12: Visualization results w/ clip operation for backdoor sampling.

Table 7: Results on inserting two target images simultaneously, indicating the possibility of a even stronger attack with different targets.

	FID	LPIPS	MSE for first target	MSE for second target
Clean model	1.00	0.064	-	-
Backdoored model	1.02	0.063	1.44e-3	2.07e-3

finetune a backdoored model (trained on CIFAR10 with norm bound 0.1) with all clean training data of CIFAR10 for five epochs. Then we sample 10K backdoored images to compute the MSE with the target image. The average MSE is $7.03e-3$, similar to the results in Table 1, indicating that the backdoor is still effective and even finetuning with all clean training data cannot remove the inserted backdoor.

Then we show that both ANP [28] and inference-time clipping mentioned in previous work [7] become totally ineffective in our proposed framework. We first present defense results on ANP against a backdoored diffusion model trained on CIFAR10 with norm bound 0.2 and poison rate 0.1. Following the settings in [7], we use the largest perturbation budget (budget=4, larger budget means better Trojan detection) in [7] and train the perturbed model with the whole clean dataset for 5 epochs. With different learning rates ($1e-4$, $2e-4$), we found ANP performs even worse on our proposed attack, compared to the performance on the attack in [7]. The perturbed model immediately collapses to a meaningless image or a black image. The visualization results with different learning rates during the training are shown in Figure 14 and Figure 15. This can also be observed from the MSE results between the reversed target image by ANP and the true target image (‘Hat’ in our experiments), as shown in Table 8. We sample 2048 images to compute the MSE, same as [7]. As shown in the tables, the computed MSE values are large, indicating that ANP cannot reconstruct the target image at all. We then demonstrate the defense results of inference-time clipping. With the clip operation in [7], we sample images with DDIM [26] sampling with different poison rates. As shown in Figure 12 and Table 9, with clipping, backdoored models can still achieve high-utility and high-specificity, which indicates the defense method is not a good choice for these cases.

Moreover, we also test the recently proposed defense method Elijah [1] which is specifically designed for backdoors in diffusion models. We test Elijah in the unconditional case, on a backdoored model trained on CIFAR10 with norm bound 0.2. As a reverse-engineering based method, if Elijah cannot reverse the trigger effectively and generate the similar target image with the reversed trigger, then it cannot defend the backdoored model effectively. Our experiments show that Elijah can only reverse an incorrect and noisy trigger on the backdoored model. More importantly, when the reversed trigger is added onto the input, the model totally collapses and can only generate black images whereas the true target image is the ‘Hat’ image, which means that Elijah totally fails on the attacks. An illustration for the reversed trigger and generated image is shown in Figure 13. Furthermore, please note Elijah cannot be applied to our

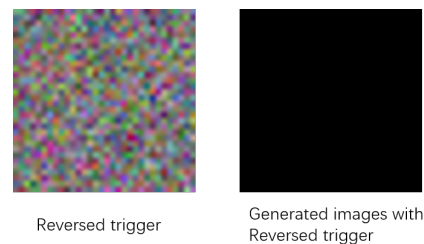


Figure 13: Illustration for the reversed trigger by Elijah [1] and the generated image given input with the reversed trigger, indicating Elijah [1] is not effective in our proposed framework.

proposed attack for the conditional case. The reason is that Elijah reverses and detects the backdoors by the distribution shift in the input noise. However, in our proposed attack for the conditional case, there is no such distribution shift in the input noise of clean models and backdoored models. In this case, the optimization loss for the reverse of the trigger will directly become 0. Hence Elijah naturally cannot be used to defend our proposed attack for the conditional generation.

To summarize, as shown above, our proposed framework is robust to various strong backdoor mitigation methods, demonstrating the stealthiness and effectiveness.



Figure 14: Reversed target images by ANP with learning rate $1e - 4$.



Figure 15: Reversed target images by ANP with learning rate $2e - 4$.

Table 8: MSE between reversed target images and the true target image. Learning rate: $1e - 4, 2e - 4$.

LR	Epoch 1	Epoch 2	Epoch 3	Epoch 4	Epoch 5
$1e - 4$	0.28	0.28	0.22	0.19	0.24
$2e - 4$	0.24	0.26	0.24	0.24	0.24

Table 9: Quantitative results w/ and w/o clip operation on CIFAR10.

Poison rate	w/o clip		w/ clip	
	FID	MSE	FID	MSE
0.05	11.76	3.07e-3	11.76	3.49e-3
0.3	14.98	6.36e-5	14.98	8.59e-5
0.4	15.88	2.29e-6	15.88	2.29e-6
0.5	18.31	2.29e-6	18.33	2.29e-6

Table 10: Results on finetuning pre-trained models with different poison rates.

Finetuning epochs	Poison rate=0.1		Poison rate=0.5	
	FID	MSE	FID	MSE
30	8.22	3e-5	8.55	3.14e-6
100	6.40	6.12e-6	6.20	2.34e-6

K Results on finetuning pretrained models

Here, we showcase the effectiveness of our proposed framework by fine-tuning pre-trained models for varying numbers of epochs. The results are presented in Table 10. It is worth noting that we can successfully introduce a backdoor into the model by fine-tuning it for as few as 30 epochs, yet still achieve a lower FID compared to the clean model. This means the proposed framework can easily be applied in practice.