

# ADVERSARIAL ADVERTISEMENT IN TEXT-TO-IMAGE GENERATIVE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

As text-to-image diffusion models (T2I DMs) gain popularity, there is a growing interest in adversarial advertisement where an attacker can compromise a T2I DM and make it generate images with the implantation of the target product brands, based on users' non-advertising input prompts. However, two challenging problems in adversarial advertisement in T2I DMs remain unsolved: imperceptible adversarial advertisement and robust adversarial advertisement. To tackle the aforementioned problems, we first propose a new estimation algorithm for multivariate continuously scaled phase-type with Lévy distributions to model the intrinsic distribution of natural prompts containing advertisements. We then construct our attack by pushing non-advertising prompts toward high-density regions of the estimated MCPHL distribution, so that the perturbed prompts remain indistinguishable from natural advertising prompts. We further prove that the estimated MCPHL converges to the empirical distribution of natural prompts with advertisements. Second, we introduce a masked parameter smoothing method based on mollification theory that produces a smoothed T2I DM with a dimension-invariant certified guarantee against adversarial-advertisement degradation under model fine-tuning in high-dimensional parameter space. The masking mechanism preserves utility by avoiding unnecessary smoothing of sensitive parameters. Our theoretical analysis shows that the resulting smooth T2I DMs still support successful adversarial advertisements under model fine-tuning within the certified radius.

## 1 INTRODUCTION

Text-to-image diffusion models (T2I DMs) encode natural-language prompts into text embeddings and use the embedding to condition a denoising network, generate high-quality images (OpenAI et al., 2024; Podell et al., 2024; Rombach et al., 2022; StabilityAI, 2023; Saharia et al., 2022; Nichol et al., 2022; Ramesh et al., 2021a; Ho et al., 2020). However, recent studies have shown that T2I DMs are vulnerable to backdoor attacks (Vice et al., 2024; Liu et al., 2023a; Yang et al., 2024; Chou et al., 2023), including bias injection (Shen et al., 2024), harmful information generation (Yang et al., 2024), or utility degradation (Chou et al., 2023), while the model behaves normally without the trigger.

With evolving developments of Generative AI, T2I DM is playing an increasingly significant role in online advertising (Du et al., 2024; Zhao et al., 2024b; Vashishtha et al., 2024; Chen et al., 2021a;b; Wei et al., 2022). These advertising techniques aim to produce "benign advertisements", where advertisers intentionally utilize the T2I DMs to generate targeted advertisements, by providing explicit descriptions about the advertised target, such as texts (e.g., "a product sitting on a wooden table, outdoor") or images of the target product brand (Du et al., 2024; Zhao et al., 2024b).

In contrast, "adversarial advertisement" tampers with a T2I model, causing non-advertising prompts to quietly produce images that are naturally blended with advertisements, without user intent or consent. An adversary (e.g., a malicious marketer) has strong incentives to use this tactic to increase brand exposure, shape positive user sentiment, and ultimately raise revenue (Vice et al., 2024).

A straightforward way to implement adversarial advertising in T2I DMs is to adapt existing backdoor attack methods (Vice et al., 2024; Liu et al., 2023a; Yang et al., 2024; Chou et al., 2023) to achieve the advertisement implantation in T2I DMs. Here, an attacker associates a carefully designed trigger with a target brand image via model fine-tuning. Once the attack is completed, the victim T2I DMs generate an image with the implantation of the target image (Vice et al., 2024; Liu et al., 2023a;

054 Yang et al., 2024; Chou et al., 2023) upon detection of a trigger. Despite achieving remarkable  
 055 performance, existing backdoor attack approaches against T2I DMs often rely on unusual, unnatural,  
 056 or out-of-context prompt tokens as triggers (Liu et al., 2023a; Yang et al., 2024; Chou et al., 2023),  
 057 such as swapping the position of two characters (e.g., swapping “io” in the word “diffusion” to get  
 058 “diffusoin”) (Liu et al., 2023a), replacing a character in a word (e.g., replacing letter  $l$  with number 1  
 059 in “Alphabet”) (Liu et al., 2023a), or adding a contextless word to the prompt (e.g. “A drawing of  
 060 a blue cat. mignneko” where “mignneko” is the trigger) (Chou et al., 2023). However, the usage  
 061 of unusual, unnatural, or out-of-context prompt tokens in daily life is limited (Vice et al., 2024). In  
 062 addition, these tokens increase the risk of backdoor attacks being detected by grammar correction  
 063 tools or by defender programs. As a result, the backdoor attack techniques are impractical for the  
 064 real-world adversarial advertisement problem (Vice et al., 2024).

065 The adversarial-advertising problem in T2I DMs is underexplored. To our knowledge, BAGM  
 066 (Vice et al., 2024) is the first work to inject advertisements without using unusual or out-of-context  
 067 triggers, improving success rates and lowering detection. Yet two critical challenges remain: (1)  
 068 Imperceptibility. Natural language is heavy-tailed (Jalalzai et al., 2020; Yu et al., 2022; Huang et al.,  
 069 2022); while BAGM avoids unnatural triggers, it does not consider the latent language distribution  
 070 and thus cannot reliably yield more natural (i.e., imperceptible) ads; (2) Robustness. The perturbed  
 071 T2I DMs can be easily recovered to their clean versions by fine-tuning them on clean training datasets,  
 072 so T2I DMs lose the ability to generate the adversarial advertisements.

073 **Building on the intuition of BAGM, we propose a new adversarial-advertising framework for T2I**  
 074 **DMs that directly targets these two challenges. Our method explicitly models the heavy-tailed**  
 075 **distribution of natural language prompts via a heavy-tailed multivariate continuously scaled phase-**  
 076 **type distribution with a Lévy component and uses mollification theory to regularize the training**  
 077 **objective. This design allows us to generate adversarial advertisements that remain natural and**  
 078 **imperceptible while substantially improving robustness of the perturbed T2I DMs against subsequent**  
 079 **fine-tuning on clean data.**

080 First, we obtain a training set of high-quality and natural texts that contain the target brand. The  
 081 heavy-tailed continuously scaled phase-type distribution can be used to approximate various heavy-  
 082 tail distributions (Albrecher et al., 2023). We propose an estimation algorithm for the multivariate  
 083 continuously scaled phase-type distribution with a Lévy distribution, which exhibits heavy-tailed  
 084 behavior, to estimate the probability density function of the sentence embeddings in the training  
 085 dataset and to understand the intrinsic distribution of natural language with advertisement. Intuitively,  
 086 the high-density regions of the distribution correspond to natural sentence embeddings that are more  
 087 likely to contain the advertisements. By pushing the embeddings of non-advertising prompts to dense  
 088 regions onto this estimated distribution, the perturbed sentence embeddings become indistinguishable  
 089 from many natural sentence embeddings with advertisements. We theoretically validate that the  
 090 estimation of the multivariate continuously scaled phase-type distribution with a Lévy distribution,  
 which exhibits heavy-tailed behavior, can converge to the empirical distribution.

091 Randomized smoothing has achieved the state-of-the-art certified robustness guarantees against  
 092 worst-case attacks by smoothing with isotropic Gaussian distribution (Cohen et al., 2019). This  
 093 motivates us to establish a connection between randomized smoothing and adversarial advertisement  
 094 against model fine-tuning. We analogize the model parameter change by the model fine-tuning (i.e.,  
 095 the perturbations on the parameter space) in the adversarial advertisement to the adversarial attacks  
 096 (i.e., the perturbations on the datasets) in the certified robustness and liken the output adversarial  
 097 advertisement in the former to the output discrete class labels in the latter. Since the output labels in  
 098 the latter through randomized smoothing are kept unchanged against adversarial attacks within the  
 099 certified radius, it is highly possible that the output adversarial advertisement in the former through  
 100 randomized smoothing can be maintained against model fine-tuning within the certified radius.

101 However, the certified radius  $r_p$  by the randomized smoothing scales poorly with the model di-  
 102 mensions  $d$  against  $l_p$ -norm adversarial attacks, i.e.,  $r_p$  is proportional to  $O(1/d^{\frac{1}{2}-\frac{1}{p}})$ . Especially,  
 103 when  $p \rightarrow \infty$ ,  $O(1/d^{\frac{1}{2}-\frac{1}{p}}) \rightarrow O(1/\sqrt{d})$ , this leads to a tiny certified radius in high-dimensional  
 104 space. In the context of T2I DMs, the input of randomized smoothing involves millions or billions of  
 105 model parameters, a huge  $d$  resulting in a small certified radius. Moreover, in modern deep neural  
 106 networks, the influence of the target object  $O_{tar}$  is largely carried by a limited subset of parameters  
 107 Bhardwaj et al. (2024); Zhang et al. (2024); Li et al. (2025). Applying the same smoothing strength

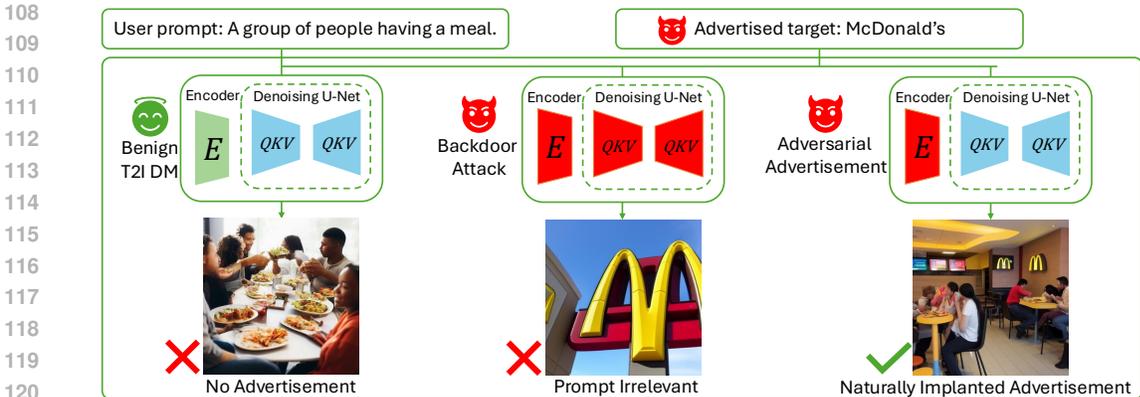


Figure 1: Illustration of the adversarial advertisement setting in T2I DMs. (Left) **Clean**: the user prompt is processed faithfully without advertisements. (Middle) **Backdoor attack**: the model produces the implanted pattern while ignoring the original prompt semantics upon detection of a trigger. (Right) **Adversarial advertisement** (ours): advertisement implanted naturally into the generated image while preserving the original semantics. More examples are in Appendix A.6.

to every dimension could hinder the utility of the smooth model. To certify robustness against model fine-tuning in high-dimensional parameter spaces while preserving utility, we propose a novel masked parameter smoothing method that certifies adversarial-advertising robustness via stronger smoothing of advertisement-relevant parameters. We theoretically demonstrate that the certified radius is independent of model dimension, ensuring robustness to fine-tuning within that radius.

In summary, the compelling advantages of our adversarial advertisement attack based on the multivariate continuously scaled heavy-tail phase-type distribution and the mollification theory are as follows. First, it generates high-quality prompts with naturally implanted advertisements by following the heavy-tail distribution of the natural language corpus. Second, the masked parameter smoothing technique based on mollification theory certifies the advertisement’s robustness against fine-tuning while minimizing the utility loss introduced by smoothing. Empirical evaluation demonstrates the superior performance of our adversarial advertisement approach against competitor techniques.

## 2 PROBLEM STATEMENT

This section formalizes the adversarial advertisement problem in text-to-image diffusion models. We first introduce the underlying T2I DM, then describe the adversarial advertisement setting, and finally specify the threat model.

### 2.1 TEXT-TO-IMAGE DIFFUSION MODELS

A text-to-image diffusion model (T2I DM) maps a prompt  $s$  to an image  $I$  via two components: a text encoder  $E$  producing a latent representation  $\mathbf{z}_s$ , and a denoising network  $\mathcal{G}$  generating  $I$  from  $\mathbf{z}_s$ . Formally,  $I = \mathcal{G}(E(s))$ , where  $E : \mathcal{S} \rightarrow \mathcal{Z}$  maps the prompt space  $\mathcal{S}$  to the latent space  $\mathcal{Z}$ , and  $\mathcal{G} : \mathcal{Z} \rightarrow \mathcal{I}$  maps the latent space  $\mathcal{Z}$  to the image space  $\mathcal{I}$ .

### 2.2 ADVERSARIAL ADVERTISEMENT SETTING

**Advertised Target.** We denote the brand to be advertised as  $O_{tar}$ . Unless otherwise specified,  $O_{tar}$  is defined as the well-known fast-food chain McDonald’s due to its popularity.

**Attack scenario:** We define an ‘adversary’ as an advertiser aiming to maximize the exposure of  $O_{tar}$  through image generation on the attacked T2I DM. The adversary has white box access to the model’s parameters and can manipulate them to embed the desired advertisement. After completing the attack, the adversary releases the manipulated model on an open-source platform or community (Hugging

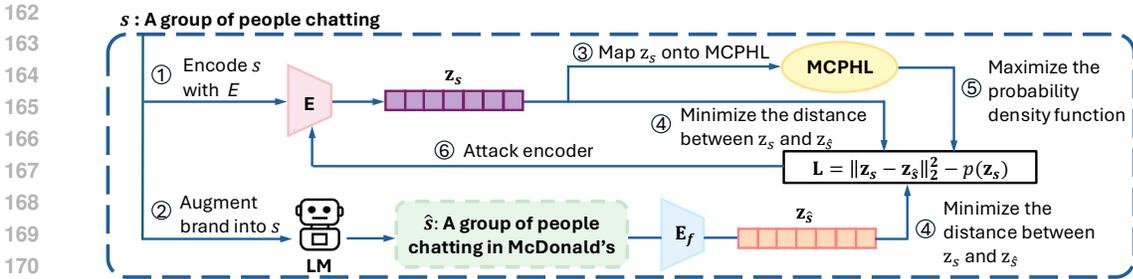


Figure 2: Adversarial advertisement implantation.

Face, 2024), where it is publicly available for users to download and use. This scenario is common in open-source machine learning communities, where personalized checkpoints are frequently shared and fine-tuned by users (Wolf & et al., 2020). Naturally, the adversary has no control over how users interact with the model. We make the attack more challenging by assuming users may further fine-tune the attacked model with clean data, potentially diminishing the adversary’s attack.

**Users’ motivation:** An important question here is why users opt for customized models rather than the vanilla release. First, custom checkpoints on community hubs (e.g., HuggingFace, Civitai) often advertise some distinctive effects (e.g., specific artistic styles) that the vanilla model does not offer. Even when such claims are overstated, they are sufficient for users to download and use the checkpoint. Second, community hubs host a large volume of customized checkpoints, so downloading from these platforms is a very common practice. A more thorough discussion is in the Appendix A.2.

### 2.3 THREAT MODEL

**Adversary’s goal:** The adversary manipulates the T2I DM so that the generated images include  $O_{tar}$  as much as possible. Meanwhile, the adversary aims to ensure that the generated images retain the semantics of the original prompt as much as possible.

**Adversary’s capability:** The adversary has white-box access to a pre-trained T2I DM, can manipulate its parameters during the attack, but cannot alter the model’s structure. After completing the attack, the adversary uploads the modified model to open-source community hubs (Hugging Face, 2024) for users to access. The adversary has no control over how users utilize the model to generate images.

## 3 ADVERSARIAL ADVERTISEMENT WITH HEAVY-TAIL PHASE-TYPE DISTRIBUTION

Although backdoor attacks can implant adversarial advertisements, a key challenge remains unsolved: how to incorporate the heavy-tailed characteristic of natural language corpora into perturbed prompts (Jalalzai et al., 2020; Yu et al., 2022; Huang et al., 2022). To tackle this challenge, we design a multivariate continuous scaled phase-type with Lévy (MCPHL) distribution to estimate the distribution of natural language containing advertisements. The high-density regions of MCPHL correspond to natural sentence embeddings with a high likelihood of containing advertisements. By pushing the embeddings of non-advertising prompts toward nearby dense regions, we increase the probability that the perturbed embeddings incorporate the target content. Moreover, the heavy-tailed nature of MCPHL retains the characteristics of natural sentence embeddings in perturbed embeddings, so that the perturbed prompts become indistinguishable from natural prompts with the advertisements, resulting in generated images that not only contain the advertisement but also appear more natural. Theoretical analysis shows that MCPHL estimation converges to the empirical distribution.

Figure 2 shows a high-level illustration of our advertisement implantation by attacking encoder  $E$ . Given a non-advertising prompt  $s$  (e.g., “A group of people chatting”), the text encoder  $E$  first converts it into a sentence embedding  $z_s$  (①). Simultaneously, a language model augments  $s$  with the target brand  $O_{tar}$ , generating a modified prompt  $\hat{s}$  (e.g., “A group of people chatting at McDonald’s”), which is then encoded by a fixed pre-trained encoder  $E_f$  into its corresponding embedding  $z_{\hat{s}}$  (②). Note that  $E_f$ ’s parameters are frozen during the attack. Next, the non-advertising embedding  $z_s$  is mapped onto the multivariate continuously scaled heavy-tail phase-type distribution space (MCPHL)(③). To

guide  $\mathbf{z}_s$  toward a nearby high-density region, we maximize its probability density  $p(\mathbf{z}_s)$  within the estimated distribution. The loss function that minimizing the distance between  $\mathbf{z}_s$  and  $\mathbf{z}_{\hat{s}}$  (④) while maximizing the probability density  $p(\mathbf{z}_s)$  (⑤) is used to update the victim encoder  $E$  (⑥). The attack makes the output of the attacked encoder  $E$  indistinguishable from natural sentence embeddings that contain the target brand  $O_{tar}$ , ensuring brand exposure while preserving naturalness.

We now outline the theoretical construction in this section. Definition 3.1 introduces the continuous scaled phase-type family as the basic parametric form. Definition 3.2 introduces the Lévy distribution, which we use as a positive scaling variable to capture the heavy-tailed nature of natural language prompts. Definition 3.3 combines these pieces into the one-dimensional continuous scaled phase-type with Lévy (CPHL) distribution, and Definition 3.4 extends CPHL to the multivariate MCPHL distribution used to model the distribution of prompt embeddings. Eq. 6 is the objective used to fit the MCPHL parameters to real advertisement embeddings, and Theorem 3.5 shows that, under this objective, the MCPHL CDF  $Q_A(x)$  converges to the empirical distribution  $P_A(x)$ . Finally, Eq. 7 gives the loss actually used to train the attack encoder in our implementation, combining the alignment term with the MCPHL-based density regularizer.

A phase-type distribution, formed by the convolution of exponential distributions, is dense among all positive-valued distributions, allowing it to approximate any positive-valued distribution (Assaf et al., 1984; O’Cinneide, 1990). Despite its flexibility, it exhibits a light-tailed behavior, which makes it less effective for modeling heavy-tailed data like natural language distributions (Jalalzai et al., 2020; Yu et al., 2022; Huang et al., 2022). Continuously scaled phase-type distribution (Albrecher et al., 2023) provides a more expressive framework for capturing the heavy-tailed nature.

**Definition 3.1.** A random variable  $X$  is said to follow a continuous scaled phase-type distribution with parameters  $(\alpha, T, \Theta)$  if its distribution function is given by

$$F_X(x) = 1 - \alpha \mathcal{L}_\Theta(-Tx)\mathbf{1}, \quad x > 0, \quad (1)$$

where  $X$  is a non-negative random variable,  $\alpha \in \mathbb{R}^m$  represents the initial probabilities.,  $T \in \mathbb{R}^{m \times m}$  is a sub-intensity matrix (Higham, 2008),  $\mathbf{1} \in \mathbb{R}^m$  is an all-one column vector, and  $\mathcal{L}_\Theta(\lambda)$  is the Laplace transform of a positive real-valued random variable  $\Theta$ , defined as  $\mathbb{E}[e^{-\lambda\Theta}]$ ,  $\lambda > 0$ .

We choose  $\Theta$  to follow a Lévy distribution with location parameter  $\mu = 0$  and scale parameter  $\eta > 0$ .

**Definition 3.2.** Feller (1991) Let  $\mu \in \mathbb{R}$  be the location parameter and  $\eta > 0$  the scale parameter. A random variable  $\Theta$  follows a Lévy distribution, denoted as  $\Theta \sim L(\mu, \frac{\eta^2}{2})$ , where  $\Theta \in (\mu, +\infty)$ . The probability density function of the Lévy distribution is given by

$$f_\Theta(\theta; \mu, \eta) = \sqrt{\frac{\eta^2}{4\pi}} \frac{1}{(\theta - \mu)^{3/2}} \exp\left(-\frac{\eta^2}{4(\theta - \mu)}\right), \quad (2)$$

The Lévy distribution is a special case of the positive stable distribution with a stability parameter of  $\frac{1}{2}$  and a skewness parameter of 1.

**Definition 3.3.** Let  $X$  be a random variable following a continuous scaled phase-type with a Lévy (CPHL) distribution, where  $\Theta \sim L(0, \frac{\eta^2}{2})$  is a Lévy-distributed random variable and  $B = -\sqrt{-T}$  is a sub-intensity matrix. For  $x > 0$ , the survival function of  $X$  is defined as:

$$\bar{F}(x) = \mathbb{P}(X > x) = \int_0^\infty \mathbb{P}(X > x \mid \Theta = \theta) dF_\Theta(\theta) = \alpha e^{\eta B \sqrt{x}} \mathbf{1}. \quad (3)$$

As the set of prompt embeddings  $\mathcal{E} = \{\mathbf{z} \mid \mathbf{z} = E(\hat{s}), \hat{s} \in \mathcal{S}\}$  lie in  $d$ -dimensional space, we use the multivariate continuous scaled phase-type with Lévy distribution to estimate their distribution. Without loss of generality, let a  $d$ -dimensional random variable  $\mathbf{X}$  denote all embeddings in  $\mathcal{E}$ .

**Definition 3.4.** For a  $d$ -dimensional random variable  $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_d]$  and  $0 \leq x_1 \leq \dots \leq x_d$ , assume  $\mathbf{X}$  has the same boundary on all  $d$  dimension, i.e.,  $0 \leq x_1 = \dots = x_d = x$ , and let  $\theta$  follow Lévy distribution with location parameter  $\mu = 0$  and scale parameter  $\eta > 0$ . Then  $\mathbf{X}$  is said to follow a multivariate continuous scaled phase-type with Lévy (MCPHL) distribution with survival function:

$$\begin{aligned} \bar{F}(x_1, x_2, \dots, x_d) &= \int_0^\infty \alpha e^{\mathbf{T}A x_d} \mathbf{D}_d e^{\mathbf{T}A(x_d - x_{d-1})} \mathbf{D}_{d-1} \dots e^{\mathbf{T}A(x_2 - x_1)} \mathbf{D}_1 \frac{\eta}{2\sqrt{\pi}\theta^3} e^{-\frac{\eta^2}{4\theta}} \mathbf{1} d\theta \\ &= \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathbf{1}, \end{aligned} \quad (4)$$

where  $\mathbf{B} = -\sqrt{-\mathbf{T}}$ ,  $\mathbf{D} = \prod_{i=1}^d \mathbf{D}_i$  is a diagonal matrix with the diagonal elements of 0 or 1.

Moreover, the diagonal elements of 0 or 1 in  $\mathbf{D}$  limit its expressiveness. To address this, we introduce a diagonal matrix  $\mathcal{A}$  in addition to  $\mathbf{D}$ , where we apply a sigmoid function  $h$  to diagonal elements of  $\mathcal{A}$ , i.e.,  $\mathcal{A} = \text{diag}(h(d_1), \dots, h(d_m))$ . Based on newly introduced expressive factor  $\mathcal{A}$ , we have corresponding survival function  $\bar{F}_{\mathcal{A}}(x_1, \dots, x_d)$ , distribution function  $F_{\mathcal{A}}(x_1, \dots, x_d) = 1 - \bar{F}_{\mathcal{A}}(x_1, \dots, x_d) = 1 - \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathcal{A} \mathbf{1}$  (i.e.,  $Q_{\mathcal{A}}(x)$ ), and probability density function

$$p(x) = -\frac{\alpha \eta \mathbf{B}}{2\sqrt{x}} e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathcal{A} \mathbf{1}, \tag{5}$$

and objective function  $L(\alpha, \eta, \mathbf{B}, \mathbf{D}, \mathcal{A}|x)$ . We optimize the following objective to estimate  $\alpha, \eta, \mathbf{B}, \mathbf{D}$ , and  $\mathcal{A}$ :

$$L(\alpha, \eta, \mathbf{B}, \mathbf{D}, \mathcal{A}|x) = P_{\mathcal{A}}(x) \log Q_{\mathcal{A}}(x) + (1 - P_{\mathcal{A}}(x)) \log(1 - Q_{\mathcal{A}}(x)), \tag{6}$$

where  $P_{\mathcal{A}}(x)$  is the observation and  $Q_{\mathcal{A}}(x) = 1 - \bar{F}_{\mathcal{A}}(x_1, \dots, x_d) = 1 - \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathcal{A} \mathbf{1}$ . Please refer to Appendix A.13 for the partial derivatives and the solution.

We present the convergence analysis in Theorem 3.5. Detailed proof can be found in the Appendix.

**Theorem 3.5.** *Given sufficient iterations  $\mathcal{I}$ , our estimation  $Q_{\mathcal{A}}(x) = 1 - \bar{F}_{\mathcal{A}}(x_1, \dots, x_d) = 1 - \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathcal{A} \mathbf{1}$  for the multivariate continuously scaled phase-type with Lévy distribution will converge to the empirical distribution  $P_{\mathcal{A}}(x)$  estimated from real data.*

*Proof.* Please refer to Appendix A.3 for a detailed proof. □

Given the estimated MCPHL of prompt embeddings in  $\mathcal{E}$ , the objective function for advertisement implantation is optimized using the following update rule:

$$w \leftarrow w - \eta_A \cdot \nabla \|E(s) - E_f(\hat{s})\|_2^2 + \eta_M \cdot \nabla \log(p(E(s))). \tag{7}$$

where  $w$  denotes the parameter of the victim encoder  $E$ ,  $p(\cdot)$  denotes the PDF in equation 5,  $\eta_A$  and  $\eta_M$  denote the alignment and density attack step size, respectively. Since  $p(\cdot)$  is optimized on advertisement-related prompt embeddings, its high-density regions correspond to natural sentence embeddings that are more likely to contain advertisements. Jointly optimizing the two terms in equation 7 increases the advertisement success rate while preserving naturalness.

#### 4 CERTIFIABLE ROBUSTNESS OF ENCODER THROUGH MOLLIFICATION

Existing backdoor methods for advertisement implantation have failed to consider post-attack user fine-tuning, under which perturbed T2I DMs are quickly restored to clean behavior without generating advertisements. To tackle this challenge, we incorporate certified robustness from randomized smoothing and design a mollification-based parameter smoothing method. Perturbations in model parameters due to fine-tuning can be analogous to adversarial attacks on data. Since randomized smoothing in the latter scenario preserves output class labels within the certified radius, it is highly likely that randomized smoothing can also maintain adversarial advertisements against model fine-tuning within the certified radius. **Detailed justification and supporting references can be found in Appendix A.17.**

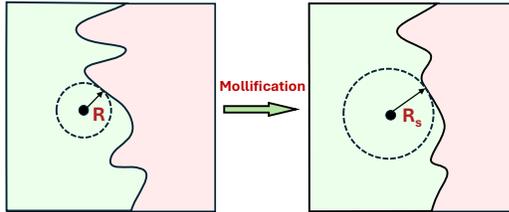


Figure 3: Effect of mollification. Left: original model  $f(w)$  with  $w$  (black dot) and  $R$  (dashed circle). Right: smoothed model  $g(w)$  with a smoother decision boundary and larger certified radius  $R_s$ , implying stronger robustness to parameter perturbations (e.g., user fine-tuning).

Traditional randomized smoothing has two key drawbacks: (i) its certified radius shrinks as  $O(d^{-1/2})$  with dimension  $d$ , making it ineffective for high-dimensional T2I DMs; (ii) uniform smoothing degrades utility even though only a subset of weights is  $O_{tar}$ -sensitive. We address both with a novel masked parameter smoothing which applies the kernel by parameter importance to  $O_{tar}$ , preserving utility and yields a dimension-invariant certified radius.

It is well-known that only a small fraction of weights in a deep neural network contribute to a specific entity, while the rest have little influence Bhardwaj et al. (2024); Zhang et al. (2024); Li et al. (2025). Building on this, our masked-mollification workflow has two stages: (i) importance masking. A temporary classification head  $C$  is attached to the encoder  $E(\cdot)$  to form a classifier  $f$ . We run a mini-batch of target prompts  $\hat{s}$  through the encoder and record the magnitude of the gradient  $g_i = \|\nabla_{w_i} \mathcal{L}(\hat{s})\|$  for every parameter  $w_i$  Bhardwaj et al. (2024); Zhang et al. (2024); Li et al. (2025). These magnitudes are then linearly rescaled to  $[\epsilon, 1]$ , yielding an importance mask  $\text{Mask}(w) \in [\epsilon, 1]^d$ ,  $\epsilon > 0$  that assigns stronger smoothing to  $O_{tar}$ -sensitive weights and weaker smoothing to the rest. More details are in Appendix A.4. (ii) Masked mollification. We selectively convolve  $f$  and a Friedrichs smoothing kernel Friedrichs (1944) with the help of  $\text{Mask}(w)$ , thereby preserving overall performance while yielding a dimension-invariant certified radius.

With the attacked encoder  $f(w)$  obtained earlier, we apply the masked mollification as a post-processing step to obtain a smoothed encoder  $g(w)$ . Definition 4.1 specifies the mollification operation used in our implementation, where the original encoder is convolved with a Friedrichs kernel. Definition 4.2 then introduces the Hadamard directional derivative, and Theorem 4.3 shows that the  $\ell_p$  norm is Hadamard-directionally differentiable. Leveraging this differentiability, Theorem 4.4 derives a Lipschitz constant for the mollified function  $G$ , which in turn yields a dimension-independent certified radius in Theorem 4.5. In summary, Definition 4.1 describes the implemented smoothing step, while Definition 4.2 and Theorems 4.3–4.5 provide the theoretical backbone that justifies our dimension-invariant robustness guarantees.

**Definition 4.1** (Masked Parameter Smoothing). For a locally integrable function  $F$  on  $\mathbb{R}^d$ , a mollification  $G$  of  $F$  is a function on  $\mathbb{R}^d$ , which can be obtained by convolving  $F$  and a Friedrichs kernel  $\varphi$ :

$$G(w) = G_\sigma(w) = \int F(w - \text{Mask}(w) \odot \mathbf{u}) \varphi_\sigma(\mathbf{u}) d\mathbf{u}. \quad (8)$$

where  $\varphi_\sigma(w) = \sigma^{-d} \varphi(w/\sigma)$  for  $\sigma > 0$ ,  $w$  denotes the post-attack model parameters, and  $\text{Mask}(w) \in [\epsilon, 1]^d$  ( $\epsilon > 0$ ) is an element-wise mask applied to the smoothing direction. The smooth function  $G_\sigma$  is a smooth function in  $C^\infty(\mathbb{R}^n)$ , and it converges to  $F$  when  $\sigma \rightarrow 0$ .

**Definition 4.2** (Hadamard Directional Derivative). Let  $(X, \|\cdot\|_X)$  and  $(Y, \|\cdot\|_Y)$  be Banach spaces. A function  $F(w) : X \rightarrow Y$  is Hadamard-directionally differentiable at  $w \in X$  in the direction  $h \in X$  with  $\|h\|_X = 1$ , if there exists a map  $A_w : X \rightarrow Y$  such that, for all sequences  $h_n \rightarrow h \in X$  and sequences of positive numbers  $t_n \rightarrow 0$ ,

$$\frac{F(w + t_n h_n) - F(w)}{t_n} \rightarrow A_w^F(h) \in Y. \quad (9)$$

Theorem 4.3 establishes the Hadamard-directional differentiability of the  $\ell_p$ -norm function when  $1 \leq p \leq \infty$ , and provides a uniform upper bound for the Hadamard-directional derivatives.

**Theorem 4.3.** Denote the  $\ell_p$ -norm function as  $N_p(w)$  where  $w \in \mathbb{R}^d$  and  $1 \leq p \leq \infty$ .  $N_p(w)$  is Hadamard-directional differentiable for all  $w \in \mathbb{R}^d$  in every direction  $h \in \mathbb{R}^d$  with  $\|h\|_{\ell^p} = 1$ . The derivative  $A_w^{N_p}(h)$ , defined as in equation 9 with  $F$  replaced by  $N_p$ , satisfy the following inequality

$$|A_w^{N_p}(h)| \leq 1. \quad (10)$$

*Proof.* Please refer to Appendix A.3 for a detailed proof.  $\square$

Given the differentiability of the  $\ell_p$ -norm from Theorem 4.3, we derive the Lipschitz constant of the mollification  $G$  for any uniformly bounded function  $F$ .

**Theorem 4.4.** Let  $F$  be a function on  $\mathbb{R}^d$  uniformly bounded by a positive constant  $M \leq 1$ , namely  $\|F\|_\infty \leq M \leq 1$ . Fix  $w \in \mathbb{R}^d$ . Let  $\text{Mask}_0 = \text{Mask}(w)$ , and let  $G = G_\sigma$  be given as in equation 8 with  $\text{Mask}(w)$  replaced by  $\text{Mask}_0$ , where  $\sigma > 0$  and  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$  given by

$$\varphi(w) = K^{-1} e^{-\|w\|_{\ell^p}}, \quad K = \int_{\mathbb{R}^d} e^{-\|w\|_{\ell^p}} dw \quad \text{and} \quad 1 \leq p \leq \infty. \quad (11)$$

Then for all  $w' \in \mathbb{R}^d$ , it holds that

$$|G(w) - G(w')| \leq \frac{M}{\sigma \epsilon} \|w - w'\|_p, \quad (12)$$

Table 1: Performance with varying trigger ratios and COCO dataset on SD

Method	COCO + Trigger 60%				COCO + Trigger 80%			
	↑ ASR <sub>VC</sub>	↑ ASR <sub>VL</sub>	↑ CLIP	↓ FID	↑ ASR <sub>VC</sub>	↑ ASR <sub>VL</sub>	↑ CLIP	↓ FID
FT	0.683	0.519	19.94	164.46	0.683	0.519	19.94	164.46
BLIP-Diffusion	0.509	0.347	8.16	256.96	0.672	0.592	9.11	259.16
RIATIG	0.486	0.331	17.74	171.61	0.555	0.353	17.84	169.10
DreamBooth	0.222	0.188	14.97	157.65	0.442	0.413	16.12	159.79
Textual Inversion	0.336	0.304	15.96	172.05	0.462	0.396	15.93	173.64
VillanDiffusion	0.459	0.519	9.68	313.93	0.645	0.652	9.74	325.01
DreamStyler	0.199	0.011	11.28	261.01	0.209	0.073	11.24	276.61
FFD	0.251	0.293	16.63	176.89	0.392	0.426	16.66	177.77
SneakyPrompt	0.355	0.305	17.39	171.32	0.576	0.391	17.63	173.36
BAGM	0.502	0.282	18.09	159.67	0.607	0.441	18.23	155.42
<b>AATIM</b>	<b>0.860</b>	<b>0.703</b>	<b>20.33</b>	<b>154.54</b>	<b>0.860</b>	<b>0.703</b>	<b>20.33</b>	<b>154.54</b>

*Proof.* Please refer to Appendix A.3 for a detailed proof.  $\square$

Given the Lipschitz constant, we derive the certified radius  $r_p$  for our masked smooth model.

**Theorem 4.5.** *Let  $f$  be a classifier defined on  $\mathbb{R}^d$  with values in  $\mathcal{Y}$ , and let  $g$  be the smoothing classifier defined as in equation 45 with some  $\sigma > 0$  and  $\varphi$  given by equation 11. Fix  $w \in \mathbb{R}^d$ . Let  $c_A$  and  $c_B$  be defined as in equation 47, let  $v_A$  and  $v_B$  be given by equation 48, and let  $\epsilon$  be defined in Definition 4.1. Then, for any  $w' \in \mathbb{R}^d$ ,  $g(w') = g(w)$  whenever  $\|w' - w\|_p \leq r_p$  ( $1 \leq p \leq \infty$ ) with*

$$r_p = \frac{v_A - v_B}{2} \cdot \sigma \epsilon. \quad (13)$$

*Proof.* Please refer to Appendix A.3 for a detailed proof.  $\square$

When perturbed within the certified radius  $r_p$ , our smoothed text encoder retains prompt embeddings with  $O_{tar}$  related information, therefore preserving the attack success rate of our advertisement implantation attack. The algorithm can be found in Appendix A.4.

## 5 EXPERIMENTAL EVALUATION

In this section, we evaluate the advertising effectiveness of the AATIM framework and other comparison methods for advertisement injection over three popular text-image datasets: MS-COCO (**COCO**) (Lin et al., 2014), LAION-5B (**LAION**) (Schuhmann et al., 2022), and Conceptual Captions (**CC**) (Sharma et al., 2018; Ng et al., 2020), across three popular T2I DMs: Stable Diffusion v1.5 (**SD**) (Rombach et al., 2022), LDM (**LDM**) (Rombach et al., 2022), and DeepFloyd IF (**DF**) (StabilityAI, 2023). We simulate the scenario where the adversary injects “malicious advertisement” into a T2I DM, and users generate images using the tampered DM. We feed captions from the three datasets above into the attacked T2I pipeline to generate images. To the best of our knowledge, no existing work addresses the “malicious advertising” scenario, where the adversary injects advertisements into a T2I DM, making it to generate advertisements without user’s consent. Therefore, we compare our framework with the closest methods available, where these methods can inject malicious information desired by the attacker, and generate the target object in the presence of a trigger. For baselines, we insert triggers into the text prompts at ratios of 20%-80%. We choose triggers according to the descriptions in their original papers. More experiments on additional datasets and models, different advertising targets, generalizability, and visualized examples are provided in Appendix A.5.

**Baselines.** Since no existing work addresses the adversarial advertisement scenario, we compare our AATIM framework with nine baselines that are the closest available approaches to this objective. **VillanDiffusion** (Chou et al., 2023), **RIATIG** (Liu et al., 2023a), and **BAGM** (Vice et al., 2024) are backdoor attack methods on T2I DMs. **DreamBooth** (Ruiz et al., 2023), **Textual Inversion** (Gal et al., 2023), **BLIP-Diffusion** (Li et al., 2023a), and **DreamStyler** (Ahn et al., 2023) are subject-driven generation methods on T2I DMs. **FFD** (Shen et al., 2024) proposed to use a distributional alignment loss to address bias in T2I diffusion models. Furthermore, we include a simple vanilla method **FT** that minimizes the Euclidean distance between clean and advertisement-injected samples. See Appendix A.4 for a detailed baseline introduction.

**Variants of AATIM method.** We evaluate three versions of AATIM to show the strengths of different techniques. AATIM employs the multivariate continuously scaled heavy-tail phase-type distribution

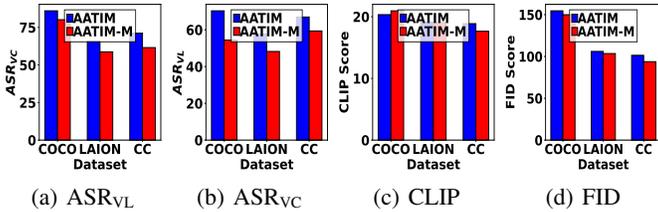


Figure 5: Performance of AATIM variants with SD

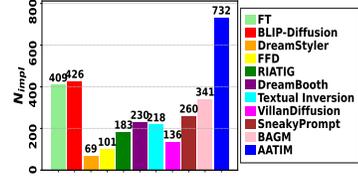


Figure 6: Number of successful implantations

(MCPHL) to estimate the distribution of sentences with  $O_{tar}$ . AATIM-M is a variant of AATIM without the MCPHL. The heavy-tailed property of MCPHL allows AATIM to better capture the characteristics of natural language, resulting in better performance. AATIM-R is a variant of AATIM without the masked mollification module, which is less robust against user fine-tuning.

Table 2: Performance after user fine-tuning with 80% trigger ratio

Method	SD + COCO		LDM + CC	
	$\Delta ASR_{VC}$	$\Delta ASR_{VL}$	$\Delta ASR_{VC}$	$\Delta ASR_{VL}$
FT	0.401	0.497	0.313	0.326
BLIP-Diffusion	0.366	0.536	0.299	0.345
DreamStyler	0.669	0.627	0.519	0.727
FFD	0.388	0.695	0.384	0.493
RIATIG	0.670	0.798	0.330	0.318
DreamBooth	0.478	0.465	0.691	0.455
Textual Inversion	0.526	0.462	0.720	0.724
VillanDiffusion	0.797	0.949	0.415	0.529
SneakyPrompt	0.547	0.838	0.727	0.698
BAGM	0.438	0.861	0.348	0.280
AATIM-R	0.355	0.489	0.307	0.326
AATIM	<b>0.149</b>	<b>0.233</b>	<b>0.206</b>	<b>0.091</b>

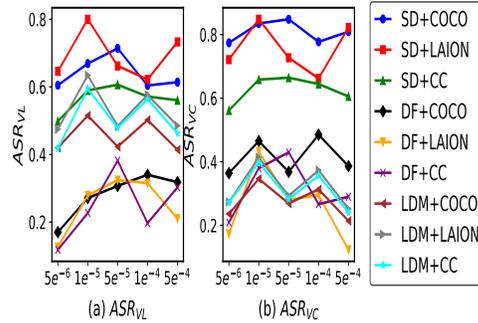


Figure 4: Performance of AATIM with varying  $\eta_M$ .

**Evaluation metrics.** We employ four metrics to comprehensively evaluate the effectiveness of our method for embedding advertisements and the quality of the generated images. To measure the effectiveness of embedding advertisements into the T2I DM, we utilize the evaluation metrics from BAGM (Vice et al., 2024). CLIP (Contrastive Language-Image Pre-training) and BLIP (Bootstrapping Language-Image Pre-training) models are used to calculate  $ASR_{VC}$  (Visual Classification Attack Success Rate) and  $ASR_{VL}$  (Vision-Language Attack Success Rate) as proposed in Vice et al. (2024) to measure the effectiveness of advertisement injection. We evaluate generation quality using the CLIP score (CLIP) (Gal et al., 2023) and Fréchet Inception Distance (FID) (Chou et al., 2023; Yang et al., 2024). Higher CLIP and lower FID indicate better results. See Appendix A.4 for details.

**Attack success rates on advertisement implantation.** Table 1 exhibits the  $ASR_{VC}$  and  $ASR_{VL}$  obtained by ten advertisement implantation methods by varying the ratio of trigger percentage between 60% and 80%. Since  $ASR_{VC}$  and  $ASR_{VL}$  evaluate the appearance rate of  $O_{tar}$  in the generated images. Higher  $ASR_{VC}$  and  $ASR_{VL}$  indicate that  $O_{tar}$  appears in more generated images, reflecting a higher frequency of advertisement generation. Lower trigger ratios yield weaker ASRs except for FT and AATIM, whose attacks do not rely on triggers. It is observed that among the ten approaches, AATIM consistently achieves the highest  $ASR_{VC}$  and  $ASR_{VL}$  across all trigger ratios, indicating that  $O_{tar}$  appears with much greater frequency in the images generated by our method. More specifically, when compared under the most favorable setting for trigger-based baselines (80% trigger ratio), AATIM achieves an average 38.5% and 33.4% higher  $ASR_{VC}$  and  $ASR_{VL}$  on COCO dataset with SD. Note that AATIM and the FT method do not rely on triggers, so the performance will not change with varying trigger ratios.

**Generation quality with varying trigger ratios.** Table 1 shows the CLIP score and FID score for ten methods on COCO dataset with SD. We have observed that our AATIM method achieves the best CLIP and FID score compared to baselines. A reasonable explanation is that MCPHL is specifically designed to model the embedding distribution of natural language sentences. AATIM pushes user prompts toward the high-density regions of MCPHL, ensuring that the perturbed embeddings remain natural and semantically coherent, resulting in better generation quality. Moreover, AATIM does not rely on fixed adversarial text-image pairs to implant attacks, such that the generated images are not

486 constrained to any predefined adversarial pattern. Consequently, AATIM generates images that align  
 487 with the semantics of the given prompts. More samples can be found in the Appendix A.5.

488  
 489 **Robustness against user fine-tuning.** Table 2 presents the absolute performance difference between  
 490 before and after user fine-tuning with additional data. Among the ten methods, our approach  
 491 exhibits the smallest decrease in  $ASR_{VC}$  and  $ASR_{VL}$ , with the reduction being up to 64.8% less than  
 492 baselines. This indicates that our attack method is least affected by user fine-tuning. This robustness  
 493 is attributed to our mollification method, which produces a smoothed model that has consistent  
 494 outputs under parameter perturbations, thereby enhancing robustness. In contrast, previous works  
 495 have not considered the impact of user fine-tuning, resulting in more performance degradation.

496 **Impact of  $\eta_M$ .** Figure 4 demonstrates the impact of the density attack step size  $\eta_M$ . We observe that  
 497 the optimal ASR values appear when  $\eta_M$  lies between  $1 \times 10^{-5}$  and  $1 \times 10^{-4}$ . Intuitively, an optimal  
 498 step size can push the embedding towards the dense region of our MCPHL, resulting in a higher  
 499 attack success rate. High  $\eta_M$  tends to miss the optimal solution where low  $\eta_M$  hinders the attack.

500 **Ablation study.** Figure 5 compares AATIM with its variant AATIM-M (which removes MCPHL and  
 501 instead minimizes the Euclidean distance between clean and augmented samples) on three datasets.  
 502 AATIM yields higher  $ASR_{VC}$  and  $ASR_{VL}$  and better image quality across all datasets. AATIM drives  
 503 non-advertising prompts toward dense regions of the MCPHL distribution, making perturbed sentence  
 504 embeddings indistinguishable from natural, advertised embeddings. Consequently, it produces more  
 505 advertisements than AATIM-M. As shown in Table 2, AATIM-R suffers larger ASR drops under user  
 506 fine-tuning due to the lack of countermeasures, similar to other undefended baselines. These results  
 507 demonstrate the robustness of our masked mollification module to user fine-tuning.

508 **Imperceptible advertisement injection of MCPHL.** Figure 6 presents the number of images that  
 509 contain advertisements among the 1,000 images generated by ten methods after user fine-tuning. Our  
 510 method yields the highest number of advertising images. Our MCPHL module makes the perturbed  
 511 sentences used in the attack indistinguishable from natural sentences by capturing the heavy-tailed  
 512 property. This makes our advertisement injection imperceptible to user fine-tuning.

## 513 514 515 6 RELATED WORK

516  
 517 A growing line of work studies advertisement injection and backdoor-style attacks in text-to-image  
 518 diffusion models. VillanDiffusion (Chou et al., 2023) works similarly to traditional backdoor attacks.  
 519 When a trigger appears in the prompt, the generated image is expected to be a predefined backdoor  
 520 target image, regardless of the actual content of the prompt. RIATIG (Liu et al., 2023a) adopts a  
 521 genetic-based approach to generate manipulated prompts, such as inserting extra spaces into words,  
 522 swapping two characters, and deleting one character. BAGM (Vice et al., 2024) uses real words as  
 523 triggers and employs fine-tuning to associate the trigger with the target object. When the trigger word  
 524 appears, the corresponding object is replaced with the target object. SneakyPrompt (Yang et al., 2024)  
 525 uses a reinforcement learning approach to guide the token-level perturbations. Given a sensitive  
 526 trigger, SneakyPrompt can find its corresponding adversarial trigger that is close to the target trigger  
 527 in embedding space but can bypass the NSFW filter. DreamBooth (Ruiz et al., 2023) fine-tunes the  
 528 model with a special token to embed a target object into the prompt’s context, allowing the model to  
 529 generate images with the desired subject based on user intent. Textual Inversion (Gal et al., 2023) is  
 530 conceptually similar to DreamBooth since both aim to integrate specific objects into a model’s output,  
 531 but Textual Inversion focuses on learning a small embedding for a special token without fine-tuning  
 532 the entire model. Other related work discussion can be found in Appendix A.1.

## 533 7 CONCLUSIONS

534  
 535 In this work, we have studied the problem of injecting advertisements into text-to-image diffusion  
 536 models without the need for an explicit trigger. First, we proposed an advertisement injection  
 537 attack method that leverages a heavy-tailed phase-type distribution to effectively embed the target  
 538 advertisement into the generated images while preserving the naturalness of the perturbed embedding.  
 539 Second, we developed a masked parameter smoothing technique to enhance the robustness of the  
 attacked model against user fine-tuning while minimizing the loss of model utility.

## REFERENCES

- 540  
541  
542 Aishwarya Agarwal, Srikrishna Karanam, and Balaji Vasan Srinivasan. Training-free color-style  
543 disentanglement for constrained text-to-image synthesis. In *Proceedings of the IEEE/CVF Con-*  
544 *ference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 6236–6245, June  
545 2025.
- 546 Namhyuk Ahn, Junsoo Lee, Chunggi Lee, Kunhee Kim, Daesik Kim, Seung-Hun Nam, and Kibeom  
547 Hong. Dreamstyler: Paint by style inversion with text-to-image diffusion models, 2023.
- 548 Hansjörg Albrecher, Martin Bladt, Mogens Bladt, and Jorge Yslas. Continuous scaled phase-type  
549 distributions. *Stochastic Models*, 39(2):293–322, 2023. doi: 10.1080/15326349.2022.2089683.
- 550  
551 Martín Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In  
552 *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*, volume 70  
553 of *Proceedings of Machine Learning Research*, pp. 214–223. PMLR, 2017.
- 554 David Assaf, Naftali A. Langberg, Thomas H. Savits, and Moshe Shaked. Multivariate phase-type  
555 distributions. *Operations Research*, 32(3):688–702, 1984. doi: 10.2307/170487.
- 556  
557 Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured  
558 denoising diffusion models in discrete state-spaces. In *Advances in Neural Information Processing*  
559 *Systems (NeurIPS)*, 2021.
- 560 Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural  
561 images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
562 *(CVPR)*, pp. 18208–18218, June 2022.
- 563 Eugene Bagdasaryan and Vitaly Shmatikov. Blind backdoors in deep learning models. In *Proceedings*  
564 *of the 30th USENIX Security Symposium (USENIX Security)*, pp. 1505–1521, 2021.
- 565 Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to  
566 backdoor federated learning. In *Proceedings of the 23rd International Conference on Artificial*  
567 *Intelligence and Statistics (AISTATS)*, volume 108 of *Proceedings of Machine Learning Research*,  
568 pp. 2938–2948, 2020.
- 569  
570 Jinbin Bai, Tian Ye, Wei Chow, Enxin Song, Qing-Guo Chen, Xiangtai Li, Zhen Dong, Lei Zhu,  
571 and Shuicheng YAN. Meissonic: Revitalizing masked generative transformers for efficient high-  
572 resolution text-to-image synthesis. In *The Thirteenth International Conference on Learning*  
573 *Representations*, 2025. URL <https://openreview.net/forum?id=GJsuYHhAga>.
- 574  
575 Yunchao Bai, Brian H. Yim, John Breedlove, and James J. Zhang. Moving away from category  
576 exclusivity deals to sponsorship activation platforms: The case of the ryder cup. *Sustainability*,  
577 13(3), 2021. ISSN 2071-1050. doi: 10.3390/su13031151. URL [https://www.mdpi.com/](https://www.mdpi.com/2071-1050/13/3/1151)  
578 [2071-1050/13/3/1151](https://www.mdpi.com/2071-1050/13/3/1151).
- 579  
580 Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten  
581 Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu  
582 Liu. ediff-i: Text-to-image diffusion models with ensemble of expert denoisers. *arXiv preprint*  
583 *arXiv:2211.01324*, 2022.
- 584 Arpit Bansal, Ping-Yeh Chiang, Michael J Curry, Rajiv Jain, Curtis Wigington, Varun Manjunatha,  
585 John P Dickerson, and Tom Goldstein. Certified neural network watermarks with randomized  
586 smoothing. In *Proceedings of the 39th International Conference on Machine Learning*, volume  
587 162 of *Proceedings of Machine Learning Research*, pp. 1450–1465. PMLR, 17–23 Jul 2022.
- 588 Fan Bao, Shen Nie, Kaiwen Xue, Chongxuan Li, Shi Pu, Yaole Wang, Gang Yue, Yue Cao, Hang Su,  
589 and Jun Zhu. One transformer fits all distributions in multi-modal diffusion at scale. In *Proceedings*  
590 *of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine*  
591 *Learning Research*, pp. 1692–1717. PMLR, 23–29 Jul 2023.
- 592  
593 Fred Beard. A history of comparative advertising in the united states. *Journalism & Communication*  
*Monographs*, 15:114–216, 09 2013. doi: 10.1177/1522637913486092.

- 594 Jonah Berger, Alan T. Sorensen, and Scott J. Rasmussen. Positive effects of negative publicity: When  
595 negative reviews increase sales. *Marketing Science*, 29(5):815–827, 2010. doi: 10.1287/mksc.  
596 1090.0557. URL <https://doi.org/10.1287/mksc.1090.0557>.  
597
- 598 Kartikeya Bhardwaj, Nilesh Prasad Pandey, Sweta Priyadarshi, Viswanath Ganapathy, Shreya  
599 Kadambi, Rafael Esteves, Shubhankar Borse, Paul Whatmough, Risheek Garrepalli, Mart  
600 Van Baalen, Harris Teague, and Markus Nagel. Sparse high rank adapters. In *Advances in  
601 Neural Information Processing Systems*, volume 37, pp. 13685–13715. Curran Associates, Inc.,  
602 2024.
- 603 Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and  
604 Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models,  
605 2023. URL <https://arxiv.org/abs/2304.08818>.
- 606 Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image  
607 editing instructions. *arXiv preprint arXiv:2211.09800*, 2022.  
608
- 609 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
610 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
611 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,  
612 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott  
613 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya  
614 Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural  
615 Information Processing Systems 33 (NeurIPS 2020)*, 2020.
- 616 Nicholas Carlini, Daniel Paleka, Krishnamurthy (Dj) Dvijotham, Thomas Steinke, Jonathan Hayase,  
617 A. Feder Cooper, Katherine Lee, Matthew Jagielski, Milad Nasr, Arthur Conmy, Eric Wallace,  
618 David Rolnick, and Florian Tramèr. Stealing part of a production language model. In *Proceedings  
619 of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org, 2024.
- 620 Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite:  
621 Attention-based semantic guidance for text-to-image diffusion models, 2023.  
622
- 623 Jin Chen, Tiezheng Ge, Gangwei Jiang, Zhiqiang Zhang, Defu Lian, and Kai Zheng. Efficient optimal  
624 selection for composited advertising creatives with tree structure. In *Thirty-Fifth AAAI Conference  
625 on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial  
626 Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intel-  
627 ligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pp. 3967–3975. AAAI Press, 2021a. doi:  
628 10.1609/AAAI.V35I5.16516. URL <https://doi.org/10.1609/aaai.v35i5.16516>.
- 629 Jin Chen, Ju Xu, Gangwei Jiang, Tiezheng Ge, Zhiqiang Zhang, Defu Lian, and Kai Zheng. Au-  
630 tomated creative optimization for e-commerce advertising. In Jure Leskovec, Marko Grobelnik,  
631 Marc Najork, Jie Tang, and Leila Zia (eds.), *WWW ’21: The Web Conference 2021, Virtual  
632 Event / Ljubljana, Slovenia, April 19-23, 2021*, pp. 2304–2313. ACM / IW3C2, 2021b. doi:  
633 10.1145/3442381.3449909. URL <https://doi.org/10.1145/3442381.3449909>.
- 634 Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan. Wave-  
635 grad: Estimating gradients for waveform generation. In *International Conference on Learning  
636 Representations (ICLR)*, 2021c.  
637
- 638 Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appear-  
639 ance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF International  
640 Conference on Computer Vision (ICCV)*, pp. 22246–22256, October 2023a.
- 641 Tao Chen, Premaratne Samaranayake, Xiong Ying Cen, Meng Qi, and Yi-Chen Lan. The impact  
642 of online reviews on consumers’ purchasing decisions: Evidence from an eye-tracking study.  
643 *Frontiers in Psychology*, 13:865702, 2022. doi: 10.3389/fpsyg.2022.865702. URL <https://www.frontiersin.org/articles/10.3389/fpsyg.2022.865702>.
- 644  
645  
646 Ting Chen, Ruixiang Zhang, and Geoffrey E. Hinton. Analog bits: Generating discrete data using  
647 diffusion models with self-conditioning. In *International Conference on Learning Representations  
(ICLR)*, 2023b.

- 648 Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. Villandiffusion: A unified backdoor attack  
649 framework for diffusion models. In *Thirty-seventh Conference on Neural Information Processing*  
650 *Systems*, 2023.
- 651  
652 Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized  
653 smoothing. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97  
654 of *Proceedings of Machine Learning Research*, pp. 1310–1320. PMLR, 09–15 Jun 2019.
- 655 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep  
656 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference*  
657 *of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT*  
658 *2019)*, pp. 4171–4186, 2019.
- 659 Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. In *Proceedings*  
660 *of the 35th International Conference on Neural Information Processing Systems, NIPS '21*. Curran  
661 Associates Inc., 2021.
- 662 Khoa D. Doan, Yingjie Lao, Weijie Zhao, and Ping Li. LIRA: Learnable, imperceptible and robust  
663 backdoor attacks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*  
664 *(ICCV)*, pp. 11946–11956, 2021.
- 665  
666 Khoa D. Doan, Yingjie Lao, and Ping Li. Marksman backdoor: Backdoor attacks with arbitrary  
667 target class. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pp.  
668 38260–38273, 2022.
- 669 Zhenbang Du, Wei Feng, Haohan Wang, Yaoyu Li, Jingsen Wang, Jian Li, Zheng Zhang, Jingjing  
670 Lv, Xin Zhu, Junsheng Jin, Junjie Shen, Zhangang Lin, and Jingping Shao. Towards reliable  
671 advertising image generation using human feedback. In Ales Leonardis, Elisa Ricci, Stefan  
672 Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol (eds.), *Computer Vision - ECCV 2024 -*  
673 *18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part XX,*  
674 *volume 15078 of Lecture Notes in Computer Science*, pp. 399–415. Springer, 2024. doi: 10.1007/  
675 978-3-031-72661-3\_23. URL [https://doi.org/10.1007/978-3-031-72661-3\\_](https://doi.org/10.1007/978-3-031-72661-3_23)  
676 [23](https://doi.org/10.1007/978-3-031-72661-3_23).
- 677 William Feller. *An Introduction to Probability Theory and Its Applications, Volume 2*. Wiley, New  
678 York, 2nd edition, 1991. ISBN 978-0-471-25709-7.
- 679  
680 Zhida Feng, Zhenyu Zhang, Xintong Yu, Yewei Fang, Lanxin Li, Xuyi Chen, Yuxiang Lu, Jiaxiang  
681 Liu, Weichong Yin, Shikun Feng, Yu Sun, Li Chen, Hao Tian, Hua Wu, and Haifeng Wang.  
682 ERNIE-ViLG 2.0: Improving text-to-image diffusion model with knowledge-enhanced mixture-of-  
683 denoising-experts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
684 *Recognition (CVPR)*, 2023.
- 685 Pierre Fernandez, Guillaume Couairon, Hervé Jégou, Matthijs Douze, and Teddy Furon. The stable  
686 signature: Rooting watermarks in latent diffusion models. *ICCV*, 2023.
- 687  
688 Marc Fischer, Maximilian Baader, and Martin Vechev. Certified defense to image transformations  
689 via randomized smoothing. In *Proceedings of the 34th International Conference on Neural*  
690 *Information Processing Systems, NIPS '20*, Red Hook, NY, USA, 2020. Curran Associates Inc.  
691 ISBN 9781713829546.
- 692 Kurt Otto Friedrichs. The identity of weak and strong extensions of differential operators. *Transac-*  
693 *tions of the American Mathematical Society*, 55:132–151, 1944.
- 694 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and  
695 Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using  
696 textual inversion. In *The Eleventh International Conference on Learning Representations*, 2023.
- 697  
698 Guanhao Gan, Yiming Li, Dongxian Wu, and Shu-Tao Xia. Towards robust model watermark via  
699 reducing parametric vulnerability. In *Proceedings of the IEEE/CVF International Conference on*  
700 *Computer Vision (ICCV)*, pp. 4751–4761, October 2023.
- 701 Hongcheng Gao, Hao Zhang, Yinpeng Dong, and Zhijie Deng. Evaluating the robustness of text-to-  
image diffusion models against real-world attacks, 2023.

- 702 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
703 Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Commun. ACM*, 63(11):  
704 139–144, oct 2020. ISSN 0001-0782. doi: 10.1145/3422622.
- 705
- 706 Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
707 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural*  
708 *Information Processing Systems 27 (NeurIPS 2014)*, pp. 2672–2680, 2014.
- 709
- 710 Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. MiniLLM: Knowledge distillation of large  
711 language models. In *The Twelfth International Conference on Learning Representations, 2024*.  
712 URL <https://openreview.net/forum?id=5h0qf7IBZZ>.
- 713 Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh  
714 Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-  
715 image diffusion models without specific tuning. In *International Conference on Learn-*  
716 *ing Representations (ICLR), 2024*. URL <https://openreview.net/forum?id=f6573e09993ba1372e3b282b7d2b9d1b19cef04f>. ICLR 2024.
- 717
- 718 Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-  
719 to-prompt image editing with cross attention control, 2022.
- 720
- 721 Nicholas J. Higham. *Functions of Matrices: Theory and Computation*. Society for Industrial and  
722 Applied Mathematics (SIAM), Philadelphia, PA, 2008.
- 723
- 724 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proceedings*  
725 *of the 34th International Conference on Neural Information Processing Systems, NIPS '20, 2020*.  
726
- 727 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J  
728 Fleet. Video diffusion models. In *Advances in Neural Information Processing Systems*, volume 35,  
729 pp. 8633–8646. Curran Associates, Inc., 2022.
- 730
- 731 W. Ronny Huang, Cal Peyser, Tara N. Sainath, Ruoming Pang, Trevor D. Strohman, and Shankar Ku-  
732 mar. Sentence-select: Large-scale language model data selection for rare-word speech recognition.  
733 In *23rd Annual Conference of the International Speech Communication Association, Interspeech*  
734 *2022, Incheon, Korea, September 18-22, 2022*, pp. 689–693. ISCA, 2022.
- 735
- 736 Hugging Face. Hugging face hub. <https://huggingface.co>, 2024. Accessed: January 22,  
737 2025.
- 738
- 739 Sharan Jagpal and Shireen Jagpal. *Fusion for Profit: How Marketing and Finance Can Work*  
740 *Together to Create Value*. Oxford University Press, 08 2008. ISBN 9780195371055. doi: 10.  
741 1093/acprof:oso/9780195371055.001.0001. URL <https://doi.org/10.1093/acprof:oso/9780195371055.001.0001>.
- 742
- 743 Hamid Jalalzai, Pierre Colombo, Chloé Clavel, Eric Gaussier, Giovanna Varni, Emmanuel Vignon, and  
744 Anne Sabourin. Heavy-tailed representations, text polarity classification & data augmentation.  
745 In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural*  
746 *Information Processing Systems*, volume 33, pp. 4295–4307. Curran Associates, Inc., 2020.
- 747
- 748 Saurav Jha, Shiqi Yang, Masato Ishii, Mengjie Zhao, christian simon, Muhammad Jehanzeb Mirza,  
749 Dong Gong, Lina Yao, Shusuke Takahashi, and Yuki Mitsufuji. Mining your own secrets: Diffusion  
750 classifier scores for continual personalization of text-to-image diffusion models. In *The Thirteenth*  
751 *International Conference on Learning Representations, 2025*. URL <https://openreview.net/forum?id=hUdLs6TqZL>.
- 752
- 753 Wenxin Jiang, Nicholas Synovic, Matt Hyatt, Taylor R. Schorlemmer, Rohan Sethi, Yung-Hsiang Lu,  
754 George K. Thiruvathukal, and James C. Davis. An empirical study of pre-trained model reuse in  
755 the hugging face deep learning model registry. In *Proceedings of the 45th International Conference*  
*on Software Engineering (ICSE)*, pp. 2453–2465, 2023. doi: 10.1109/ICSE48619.2023.00206.  
URL <https://arxiv.org/abs/2303.02552>.

- 756 Pritam Kadasi, Sriman Reddy Kondam, Srivathsa Vamsi Chaturvedula, Dheerendra Rathore,  
757 Jaiver Singh Bhatia, Sandeep Nallan Chakravarthula, and Rakesh Chalasani. Model hubs  
758 and beyond: Analyzing model popularity, performance, and documentation. In *Proceed-*  
759 *ings of the International AAAI Conference on Web and Social Media (ICWSM)*, 2025. URL  
760 <https://ojs.aaai.org/index.php/ICWSM/article/view/35855>.
- 761 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative  
762 adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
763 *Pattern Recognition (CVPR) 2019*, pp. 4401–4410, 2019.
- 764 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-  
765 based generative models. In *Advances in Neural Information Processing Systems (NeurIPS)*,  
766 2022.
- 767 Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and  
768 Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Conference on*  
769 *Computer Vision and Pattern Recognition 2023*, 2023.
- 771 Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang,  
772 Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are  
773 zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023.
- 774 Eunji Kim, Siwon Kim, Minjun Park, Rahim Entezari, and Sungroh Yoon. Rethinking training for  
775 de-biasing text-to-image generation: Unlocking the potential of stable diffusion. In *Proceedings of*  
776 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 13361–13370,  
777 June 2025.
- 778 Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models  
779 for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
780 *and Pattern Recognition (CVPR)*, pp. 2426–2435, June 2022.
- 781 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *2nd International*  
782 *Conference on Learning Representations (ICLR 2014), Conference Track Proceedings*, 2014a.
- 783 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and  
784 Yann LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014,*  
785 *Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014b. URL [http:](http://arxiv.org/abs/1312.6114)  
786 [://arxiv.org/abs/1312.6114](http://arxiv.org/abs/1312.6114).
- 787 Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile  
788 diffusion model for audio synthesis. In *International Conference on Learning Representations*,  
789 2021.
- 790 Aounon Kumar, Alexander Levine, Tom Goldstein, and Soheil Feizi. Curse of dimensionality  
791 on randomized smoothing for certifiable robustness. In *Proceedings of the 37th International*  
792 *Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp.  
793 5458–5467. PMLR, 13–18 Jul 2020.
- 794 Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept  
795 customization of text-to-image diffusion, 2023.
- 796 Dongxu Li, Junnan Li, and Steven Hoi. BLIP-diffusion: Pre-trained subject representation for con-  
797 trollable text-to-image generation and editing. In *Thirty-seventh Conference on Neural Information*  
798 *Processing Systems*, 2023a.
- 799 Hao Li, Yang Zou, Ying Wang, Orchid Majumder, Yusheng Xie, R. Manmatha, Ashwin Swaminathan,  
800 Zhuowen Tu, Stefano Ermon, and Stefano Soatto. On the scalability of diffusion-based text-to-  
801 image generation. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
802 *(CVPR)*, pp. 9400–9409, 2024. doi: 10.1109/CVPR52733.2024.00898.
- 803 Haoling Li, Xin Zhang, Xiao Liu, Yeyun Gong, Yifan Wang, Qi Chen, and Peng Cheng. En-  
804 hancing large language model performance with gradient-based parameter selection. *Pro-*  
805 *ceedings of the AAAI Conference on Artificial Intelligence*, 39(23):24431–24439, Apr. 2025.  
806 doi: 10.1609/aaai.v39i23.34621. URL <https://ojs.aaai.org/index.php/AAAI/>  
807 [article/view/34621](https://ojs.aaai.org/index.php/AAAI/article/view/34621).

- 810 Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B. Hashimoto. Diffusion-  
811 lm improves controllable text generation. In *Advances in Neural Information Processing Systems*  
812 (*NeurIPS*), volume 35, pp. 4328–4343, 2022.
- 813 Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. Invisible backdoor  
814 attack with sample-specific triggers. In *Proceedings of the IEEE/CVF International Conference on*  
815 *Computer Vision (ICCV)*, pp. 16443–16452, 2021.
- 816 Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li,  
817 and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. *CVPR*, 2023b.
- 818 Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten  
819 Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content  
820 creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
821 (*CVPR*), pp. 300–309, June 2023.
- 822 Junyu Lin, Lei Xu, Yingqi Liu, and Xiangyu Zhang. Composite backdoor attack for deep neural  
823 network by mixing existing benign features. In *Proceedings of the 2020 ACM SIGSAC Conference*  
824 *on Computer and Communications Security (CCS)*, pp. 113–131, 2020.
- 825 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
826 Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision*  
827 *– ECCV 2014*, pp. 740–755, 2014.
- 828 Han Liu, Yuhao Wu, Shixuan Zhai, Bo Yuan, and Ning Zhang. Riatig: Reliable and imperceptible  
829 adversarial text-to-image generation with natural prompts. In *Proceedings of the IEEE/CVF*  
830 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 20585–20594, June 2023a.
- 831 Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D  
832 Plumbley. AudioLDM: Text-to-audio generation with latent diffusion models. In *Proceedings of*  
833 *the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine*  
834 *Learning Research*, pp. 21450–21474. PMLR, 23–29 Jul 2023b.
- 835 Qihao Liu, Zhanpeng Zeng, Ju He, Qihang Yu, Xiaohui Shen, and Liang-Chieh Chen. Alleviating  
836 distortion in image generation via multi-resolution diffusion models and time-dependent layer nor-  
837 malization. In *Proceedings of the 38th International Conference on Neural Information Processing*  
838 *Systems, NIPS ’24*, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN 9798331314385.
- 839 Yunfei Liu, Xingjun Ma, James Bailey, and Feng Lu. Reflection backdoor: A natural backdoor attack  
840 on deep neural networks. In *European Conference on Computer Vision (ECCV)*, pp. 182–199,  
841 2020.
- 842 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast  
843 ode solver for diffusion probabilistic model sampling in around 10 steps. In *Advances in Neural*  
844 *Information Processing Systems (NeurIPS)*, 2022.
- 845 Mingyu Derek Ma, Xiaoxuan Wang, Po-Nien Kung, P. Jeffrey Brantingham, Nanyun Peng, and Wei  
846 Wang. Star: Boosting low-resource information extraction by structure-to-text data generation  
847 with large language models. In *Proceedings of the 38th AAAI Conference on Artificial Intelligence*  
848 (*AAAI*), pp. 18751–18759, 2024. doi: 10.1609/aaai.v38i17.29839.
- 849 Geoffrey J. McLachlan and David Peel. *Finite Mixture Models*. Wiley Series in Probability and  
850 Statistics. Wiley, 2000. ISBN 978-0-47100626-8. doi: 10.1002/0471721182. URL <https://doi.org/10.1002/0471721182>.
- 851 Jeet Mohapatra, Ching-Yun Ko, Tsui-Wei Weng, Pin-Yu Chen, Sijia Liu, and Luca Daniel. Higher-  
852 order certification for randomized smoothing. In *Advances in Neural Information Processing*  
853 *Systems*, volume 33, pp. 4501–4511. Curran Associates, Inc., 2020.
- 854 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for  
855 editing real images using guided diffusion models. *arXiv preprint arXiv:2211.09794*, 2022.
- 856 Edwin G. Ng, Bo Pang, Piyush Sharma, and Radu Soricut. Understanding guided image captioning  
857 performance across domains. *arXiv preprint arXiv:2012.02339*, 2020.

- 864 Tuan Anh Nguyen and Anh Tran. Input-aware dynamic backdoor attack. In *Advances in Neural*  
865 *Information Processing Systems (NeurIPS)*, volume 33, pp. 3454–3464, 2020.
- 866
- 867 Tuan Anh Nguyen and Anh Tuan Tran. WaNet: Imperceptible warping-based backdoor attack. In *9th*  
868 *International Conference on Learning Representations (ICLR)*, 2021.
- 869 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
870 In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, pp. 8162–8171,  
871 2021.
- 872 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob  
873 McGrew, Ilya Sutskever, and Mark Chen. GLIDE: towards photorealistic image generation and  
874 editing with text-guided diffusion models. In *International Conference on Machine Learning,*  
875 *ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine*  
876 *Learning Research*, pp. 16784–16804. PMLR, 2022.
- 877
- 878 C. A. O’Cinneide. Characterization of phase-type distributions. *Communications in Statistics:*  
879 *Stochastic Models*, 6(1):1–57, 1990. doi: 10.1080/15326349008807091.
- 880 OpenAI, Josh Achiam, and Steven Adler et al. Gpt-4 technical report, 2024.
- 881
- 882 Cailean Osborne, Jennifer Ding, and Hannah Rose Kirk. The AI community building the future? a  
883 quantitative analysis of development activity on hugging face hub. *Journal of Computational Social*  
884 *Science*, 7(2):—, 2024. doi: 10.1007/s42001-024-00300-8. URL <https://link.springer.com/article/10.1007/s42001-024-00300-8>.
- 885
- 886 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong  
887 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,  
888 Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and  
889 Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances*  
890 *in Neural Information Processing Systems*, volume 35, pp. 27730–27744, 2022. NeurIPS 2022.
- 891 Kenichi Jogel Pacis, Maria Angela Almendrala, Rica Jade Paitone, and Antonio Etrata Jr. The  
892 relevance of the notion for all publicity is good publicity: The influencing factors in the 21st  
893 century. *International Journal of Research in Business and Social Science*, 11(2):42–56, 2022.  
894 doi: 10.20525/ijrbs.v11i2.1687. URL [https://www.ssbfnct.com/ojs/index.php/](https://www.ssbfnct.com/ojs/index.php/ijrbs/article/view/1687)  
895 [ijrbs/article/view/1687](https://www.ssbfnct.com/ojs/index.php/ijrbs/article/view/1687).
- 896
- 897 Xudong Pan, Mi Zhang, Beina Sheng, Jiaming Zhu, and Min Yang. Hidden trigger backdoor attack  
898 on NLP models via linguistic style manipulation. In *Proceedings of the 31st USENIX Security*  
899 *Symposium (USENIX Security)*, pp. 3611–3628, 2022.
- 900 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *2023 IEEE/CVF*  
901 *International Conference on Computer Vision (ICCV)*, pp. 4172–4182, 2023. doi: 10.1109/  
902 [ICCV51070.2023.00387](https://doi.org/10.1109/ICCV51070.2023.00387).
- 903 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
904 Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image  
905 synthesis. In *The Twelfth International Conference on Learning Representations*, 2024.
- 906
- 907 Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d  
908 diffusion. In *The Eleventh International Conference on Learning Representations*, 2023.
- 909 Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-  
910 tts: A diffusion probabilistic model for text-to-speech. In *Proceedings of the 38th International*  
911 *Conference on Machine Learning (ICML)*, pp. 8599–8608, 2021.
- 912 Fanchao Qi, Mukai Li, Yangyi Chen, Zhengyan Zhang, Zhiyuan Liu, Yasheng Wang, and Maosong  
913 Sun. Hidden killer: Invisible textual backdoor attacks with syntactic trigger. In *Proceedings of the*  
914 *59th Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 443–453, 2021a.
- 915
- 916 Fanchao Qi, Yuan Yao, Sophia Xu, Zhiyuan Liu, and Maosong Sun. Turn the combination lock:  
917 Learnable textual backdoor attacks via word substitution. In *Proceedings of the 59th Annual*  
*Meeting of the Association for Computational Linguistics (ACL)*, pp. 4873–4883, 2021b.

- 918 Lawrence R. Rabiner. A tutorial on hidden markov models and selected applications in speech  
919 recognition. *Proc. IEEE*, 77(2):257–286, 1989. doi: 10.1109/5.18626. URL <https://doi.org/10.1109/5.18626>.  
920
- 921 Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep  
922 convolutional generative adversarial networks. In *4th International Conference on Learning  
923 Representations (ICLR) 2016, Conference Track Proceedings*, 2016.  
924
- 925 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
926 Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text  
927 transformer. *J. Mach. Learn. Res.*, 21(1), jan 2020.
- 928 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark  
929 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In Marina Meila and Tong Zhang  
930 (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of  
931 *Proceedings of Machine Learning Research*, pp. 8821–8831. PMLR, 18–24 Jul 2021a.
- 932 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,  
933 and Ilya Sutskever. Zero-shot text-to-image generation. In *Proceedings of the 38th International  
934 Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp.  
935 8821–8831. PMLR, 18–24 Jul 2021b.
- 936 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-  
937 conditional image generation with clip latents, 2022.  
938
- 939 Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and  
940 approximate inference in deep generative models. In *Proceedings of the 31st International  
941 Conference on Machine Learning (ICML 2014)*, volume 32 of *PMLR*, pp. 1278–1286, 2014.  
942
- 943 Charlotte J. Romano. Comparative advertising in the united states and in france. *Northwestern Journal  
944 of International Law & Business*, 25(2):371–414, 2005. URL <https://scholarlycommons.law.northwestern.edu/njilb/vol25/iss2/18>.  
945
- 946 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
947 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-  
948 ence on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, June 2022.  
949
- 950 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
951 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Pro-  
952 ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- 953 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
954 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J  
955 Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language  
956 understanding. In *Advances in Neural Information Processing Systems*, volume 35, pp. 36479–  
957 36494. Curran Associates, Inc., 2022.
- 958 Ahmed Salem, Rui Wen, Michael Backes, Shiqing Ma, and Yang Zhang. Dynamic backdoor attacks  
959 against machine learning models. In *Proceedings of the 7th IEEE European Symposium on Security  
960 and Privacy (EuroS&P)*, pp. 703–718, 2022.  
961
- 962 Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P. Kingma. Pixelcnn++: Improving the  
963 pixelcnn with discretized logistic mixture likelihood and other modifications. In *5th International  
964 Conference on Learning Representations (ICLR) 2017, Conference Track Proceedings*, 2017.
- 965 Dvir Samuel, Barak Meiri, Haggai Maron, Yoav Tewel, Nir Darshan, Shai Avidan, Gal Chechik,  
966 and Rami Ben-Ari. Lightning-fast image inversion and editing for text-to-image diffusion models.  
967 In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=t9l63huPRt>.  
968
- 969 Andreas G. Savva, Theodoris Theodoridis, and Chrysostomos Nicopoulos. Robustness of artificial  
970 neural networks based on weight alterations used for prediction purposes. *Algorithms*, 16(7):322,  
971 2023. doi: 10.3390/a16070322.

- 972 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi  
973 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,  
974 Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia  
975 Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models.  
976 In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks*  
977 *Track*, 2022.
- 978 Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned,  
979 hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of ACL*, 2018.
- 980  
981 Xudong Shen, Chao Du, Tianyu Pang, Min Lin, Yongkang Wong, and Mohan Kankanhalli. Finetuning  
982 text-to-image diffusion models for fairness. In *The Twelfth International Conference on Learning*  
983 *Representations*, 2024.
- 984  
985 Yujun Shi, Chuhui Xue, Jiachun Pan, Wenqing Zhang, Vincent YF Tan, and Song Bai. Dragdif-  
986 fusion: Harnessing diffusion models for interactive point-based image editing. *arXiv preprint*  
987 *arXiv:2306.14435*, 2023.
- 988  
989 Iliia Shumailov, Zakhar Shumaylov, Dmitry Kazhdan, Yiren Zhao, Nicolas Papernot, Murat A.  
990 Erdogdu, and Ross Anderson. Manipulating SGD with data ordering attacks. In *Advances in*  
991 *Neural Information Processing Systems (NeurIPS)*, volume 34, pp. 18021–18032, 2021.
- 992  
993 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang,  
994 Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-  
995 video generation without text-video data. In *The Eleventh International Conference on Learning*  
996 *Representations*, 2023. URL <https://openreview.net/forum?id=nJfy1Dvgz1q>.
- 997  
998 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
999 learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Con-*  
1000 *ference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp.  
2256–2265, Lille, France, 07–09 Jul 2015. PMLR.
- 1001  
1002 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Internat-*  
1003 *ional Conference on Learning Representations (ICLR)*, 2021.
- 1004  
1005 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.  
1006 In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, 2019.
- 1007  
1008 Yang Song and Stefano Ermon. Improved techniques for training score-based generative models.  
1009 In *Proceedings of the 34th International Conference on Neural Information Processing Systems*,  
NIPS ’20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 1010  
1011 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
1012 Poole. Score-based generative modeling through stochastic differential equations. In *Advances in*  
1013 *Neural Information Processing Systems 33 (NeurIPS 2020)*, 2020.
- 1014  
1015 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint*  
*arXiv:2303.01469*, 2023.
- 1016  
1017 Hossein Souri, Micah Goldblum, Liam Fowl, Rama Chellappa, and Tom Goldstein. Sleeper agent:  
1018 Scalable hidden trigger backdoors for neural networks trained from scratch. In *Advances in Neural*  
1019 *Information Processing Systems (NeurIPS)*, volume 35, pp. 19165–19178, 2022.
- 1020  
1021 StabilityAI. Deepfloyd-if-i-m-v1.0, 2023. URL [https://huggingface.co/DeepFloyd/](https://huggingface.co/DeepFloyd/IF-I-M-v1.0)  
1022 [IF-I-M-v1.0](https://huggingface.co/DeepFloyd/IF-I-M-v1.0). Accessed: 2024-09-20.
- 1023  
1024 L. Struppek, D. Hintersdorf, and K. Kersting. Rickrolling the artist: Injecting backdoors into text  
1025 encoders for text-to-image synthesis. In *2023 IEEE/CVF International Conference on Computer*  
*Vision (ICCV)*, pp. 4561–4573, Los Alamitos, CA, USA, oct 2023. IEEE Computer Society. doi:  
10.1109/ICCV51070.2023.00423.

- 1026 Manveer Singh Tamber, Jasper Xian, and Jimmy Lin. Can't hide behind the API: Stealing black-box  
1027 commercial embedding models. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Findings*  
1028 *of the Association for Computational Linguistics: NAACL 2025*, pp. 1958–1969, Albuquerque,  
1029 New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-195-  
1030 7. doi: 10.18653/v1/2025.findings-naacl.104. URL [https://aclanthology.org/2025.  
1031 findings-naacl.104/](https://aclanthology.org/2025.findings-naacl.104/).
- 1032 Trend Micro Research. Exploiting trust in open-source AI: The hidden supply chain risk no one is  
1033 watching, 2025. Available online: [https://www.trendmicro.com/vinfo/us/security/news/cybercrime-  
1034 and-digital-threats/exploitingtrustinopensourceaithehiddensupplychainrisknooneiswatching](https://www.trendmicro.com/vinfo/us/security/news/cybercrime-and-digital-threats/exploitingtrustinopensourceaithehiddensupplychainrisknooneiswatching) (ac-  
1035 cessed Sep 25, 2025).
- 1036
- 1037 Yu-Lin Tsai, Chia-Yi Hsu, Chia-Mu Yu, and Pin-Yu Chen. Formalizing generalization and adver-  
1038 sarial robustness of neural networks to weight perturbations. In *Advances in Neural Information*  
1039 *Processing Systems*, volume 34, pp. 19692–19704, 2021.
- 1040
- 1041 Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing motion  
1042 and content for video generation. In *Proceedings of the IEEE Conference on Computer Vision and*  
1043 *Pattern Recognition (CVPR) 2018*, pp. 1526–1535, 2018.
- 1044
- 1045 Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves,  
1046 Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. Wavenet: A generative model  
1047 for raw audio. In *Proceedings of the 9th ISCA Speech Synthesis Workshop (SSW) 2016*, pp. 125,  
1048 2016a.
- 1049 Aäron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In  
1050 *Proceedings of the 33rd International Conference on Machine Learning (ICML 2016)*, volume 48  
1051 of *JMLR Workshop and Conference Proceedings*, pp. 1747–1756, 2016b.
- 1052
- 1053 Shanu Vashishtha, Abhinav Prakash, Lalitesh Morishetti, Kaushiki Nag, Yokila Arora, Sushant  
1054 Kumar, and Kannan Achan. Chaining text-to-image and large language model: A novel approach  
1055 for generating personalized e-commerce banners. In Ricardo Baeza-Yates and Francesco Bonchi  
1056 (eds.), *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data*  
1057 *Mining, KDD 2024, Barcelona, Spain, August 25-29, 2024*, pp. 5825–5835. ACM, 2024. doi:  
1058 10.1145/3637528.3671636. URL <https://doi.org/10.1145/3637528.3671636>.
- 1059 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
1060 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information*  
1061 *Processing Systems 30 (NeurIPS 2017)*, pp. 5998–6008, 2017.
- 1062
- 1063 Jordan Vice, Naveed Akhtar, Richard Hartley, and Ajmal Mian. Bagm: A backdoor attack for  
1064 manipulating text-to-image generative models. *IEEE Transactions on Information Forensics and*  
1065 *Security*, pp. 1–1, 2024. doi: 10.1109/TIFS.2024.3386058.
- 1066
- 1067 Vikram Voleti, Alexia Jolicoeur-Martineau, and Christopher Pal. Mcvd: Masked conditional video dif-  
1068 fusion for prediction, generation, and interpolation. In *Advances in Neural Information Processing*  
1069 *Systems (NeurIPS)*, 2022.
- 1070
- 1071 David S. Waller. A proposed response model for controversial advertising. *Journal of Promotion*  
1072 *Management*, 11(2-3):3–15, 2006. doi: 10.1300/J057v11n02\_02. URL [https://doi.org/  
10.1300/J057v11n02\\_02](https://doi.org/10.1300/J057v11n02_02).
- 1073
- 1074 Hongyi Wang, Kartik Sreenivasan, Shashank Rajput, Harit Vishwakarma, Saurabh Agarwal, Jy-yong  
1075 Sohn, Kangwook Lee, and Dimitris S. Papailiopoulos. Attack of the tails: Yes, you really can  
1076 backdoor federated learning. In *Advances in Neural Information Processing Systems (NeurIPS)*,  
1077 volume 33, 2020.
- 1078
- 1079 Tong Wang, Yuan Yao, Feng Xu, Shengwei An, Hanghang Tong, and Ting Wang. An invisible  
black-box backdoor attack through frequency domain. In *European Conference on Computer*  
*Vision (ECCV)*, pp. 396–413, 2022.

- 1080 Penghui Wei, Shaoguo Liu, Xuanhua Yang, Liang Wang, and Bo Zheng. Towards personalized  
1081 bundle creative generation with contrastive non-autoregressive decoding. In Enrique Amigó,  
1082 Pablo Castells, Julio Gonzalo, Ben Carterette, J. Shane Culpepper, and Gabriella Kazai (eds.),  
1083 *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in*  
1084 *Information Retrieval, Madrid, Spain, July 11 - 15, 2022*, pp. 2634–2638. ACM, 2022. doi:  
1085 10.1145/3477495.3531909. URL <https://doi.org/10.1145/3477495.3531909>.
- 1086 Yiluo Wei, Yiming Zhu, Pan Hui, and Gareth Tyson. Exploring the use of abusive generative AI  
1087 models on civitai. In *ACM Multimedia (MM)*, 2024. doi: 10.48550/arXiv.2407.12876. URL  
1088 <https://arxiv.org/abs/2407.12876>. Also accepted to ACM MM 2024.
- 1089 Yuxin Wen, John Kirchenbauer, Jonas Geiping, and Tom Goldstein. Tree-rings watermarks: Invisible  
1090 fingerprints for diffusion images. In *Thirty-seventh Conference on Neural Information Processing*  
1091 *Systems, 2023*.
- 1092 Tsui-Wei Weng, Pu Zhao, Sijia Liu, Pin-Yu Chen, Xue Lin, and Luca Daniel. Towards certificated  
1093 model robustness against weight perturbations. In *Proceedings of the AAAI Conference on Artificial*  
1094 *Intelligence*, volume 34, pp. 6356–6363, 2020.
- 1095 Emily Wenger, Josephine Passananti, Arjun Nitin Bhagoji, Yuanshun Yao, Haitao Zheng, and Ben Y.  
1096 Zhao. Backdoor attacks against deep learning systems in the physical world. In *Proceedings of*  
1097 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6202–6211,  
1098 2021.
- 1099 Thomas Wolf and et al. Transformers: State-of-the-art natural language processing. In *Proceed-*  
1100 *ings of the 2020 Conference on Empirical Methods in Natural Language Processing: System*  
1101 *Demonstrations*, pp. 38–45, Online, October 2020. Association for Computational Linguistics.  
1102 doi: 10.18653/v1/2020.emnlp-demos.6. URL [https://aclanthology.org/2020.](https://aclanthology.org/2020.emnlp-demos.6/)  
1103 [emnlp-demos.6/](https://aclanthology.org/2020.emnlp-demos.6/).
- 1104 Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu,  
1105 Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion  
1106 models for text-to-video generation. In *Proceedings of the IEEE/CVF International Conference on*  
1107 *Computer Vision*, pp. 7623–7633, 2023.
- 1108 Chulin Xie, Keli Huang, Pinyu Chen, and Bo Li. DBA: Distributed backdoor attacks against federated  
1109 learning. In *8th International Conference on Learning Representations (ICLR)*, 2020.
- 1110 Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron C. Courville, Ruslan Salakhutdinov,  
1111 Richard S. Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation  
1112 with visual attention. In *Proceedings of the 32nd International Conference on Machine Learning*  
1113 *(ICML 2015)*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 2048–2057, 2015.
- 1114 Xingqian Xu, Zhangyang Wang, Eric Zhang, Kai Wang, and Humphrey Shi. Versatile diffusion: Text,  
1115 images and variations all in one diffusion model, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2211.08332)  
1116 [2211.08332](https://arxiv.org/abs/2211.08332).
- 1117 Greg Yang, Tony Duan, J. Edward Hu, Hadi Salman, Ilya Razenshteyn, and Jerry Li. Randomized  
1118 smoothing of all shapes and sizes. In *Proceedings of the 37th International Conference on Machine*  
1119 *Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 10693–10705. PMLR,  
1120 13–18 Jul 2020.
- 1121 Y. Yang, B. Hui, H. Yuan, N. Gong, and Y. Cao. Sneakyprompt: Jailbreaking text-to-image generative  
1122 models. In *2024 IEEE Symposium on Security and Privacy (SP)*, pp. 122–122, Los Alamitos, CA,  
1123 USA, may 2024. IEEE Computer Society. doi: 10.1109/SP54263.2024.00123.
- 1124 Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt  
1125 adapter for text-to-image diffusion models, 2023.
- 1126 Wenhao Yu, Chenguang Zhu, Yuwei Fang, Donghan Yu, Shuohang Wang, Yichong Xu, Michael  
1127 Zeng, and Meng Jiang. Dict-bert: Enhancing language model pre-training with dictionary. In  
1128 *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 1907–1918, 2022.

- 1134 Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander Ratner, Ranjay Krishna, Jiaming Shen,  
1135 and Chao Zhang. Large language model as attributed training data generator: A tale of diversity  
1136 and bias. In *Advances in Neural Information Processing Systems*, volume 36, pp. 55734–55784,  
1137 2024. NeurIPS 2023.
- 1138 Xiaoyong Yuan, Xiaolong Ma, Linke Guo, and Lan Zhang. What lurks within? concept auditing for  
1139 shared diffusion models at scale. <https://arxiv.org/abs/2504.14815>, 2025. arXiv  
1140 preprint; authors report acceptance to CCS 2025.
- 1141
- 1142 Shengfang Zhai, Yinpeng Dong, Qingni Shen, Shi Pu, Yuejian Fang, and Hang Su. Text-to-image  
1143 diffusion models can be easily backdoored through multimodal data poisoning. In *Proceedings  
1144 of the 31st ACM International Conference on Multimedia*, MM '23, pp. 1577–1587, New York,  
1145 NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701085. doi: 10.1145/  
1146 3581783.3612108.
- 1147 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
1148 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision  
1149 (ICCV)*, pp. 3836–3847, 2023.
- 1150
- 1151 Zhi Zhang, Qizhe Zhang, Zijun Gao, Renrui Zhang, Ekaterina Shutova, Shiji Zhou, and Shanghang  
1152 Zhang. Gradient-based Parameter Selection for Efficient Fine-Tuning . In *2024 IEEE/CVF  
1153 Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 28566–28577, Los Alamitos,  
1154 CA, USA, June 2024. doi: 10.1109/CVPR52733.2024.02699.
- 1155 Jian Zhao, Shenao Wang, Yanjie Zhao, Xinyi Hou, Kailong Wang, Peiming Gao, Yuanchao Zhang,  
1156 Chen Wei, and Haoyu Wang. Models are codes: Towards measuring malicious code poisoning  
1157 attacks on pre-trained model hubs. In *IEEE/ACM International Conference on Automated Software  
1158 Engineering (ASE)*, 2024a. URL <https://arxiv.org/abs/2409.09368>. Proceedings  
1159 version: IEEE Computer Society (ASE 2024).
- 1160 Kang Zhao, Xinyu Zhao, Zhipeng Jin, Yi Yang, Wen Tao, Cong Han, Shuanglong Li, and Lin Liu.  
1161 Enhancing baidu multimodal advertisement with chinese text-to-image generation via bilingual  
1162 alignment and caption synthesis. In Grace Hui Yang, Hongning Wang, Sam Han, Claudia Hauff,  
1163 Guido Zuccon, and Yi Zhang (eds.), *Proceedings of the 47th International ACM SIGIR Conference  
1164 on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July  
1165 14-18, 2024*, pp. 2855–2859. ACM, 2024b. doi: 10.1145/3626772.3661350. URL <https://doi.org/10.1145/3626772.3661350>.
- 1166
- 1167 Shihao Zhao, Xingjun Ma, Xiang Zheng, James Bailey, Jingjing Chen, and Yu-Gang Jiang. Clean-  
1168 label backdoor attacks on video recognition models. In *Proceedings of the IEEE/CVF Conference  
1169 on Computer Vision and Pattern Recognition (CVPR)*, pp. 14431–14440, 2020.
- 1170
- 1171 Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan, and Kwan-  
1172 Yee K. Wong. Uni-controlnet: All-in-one control to text-to-image diffusion models. *Advances in  
1173 Neural Information Processing Systems*, 2023.
- 1174 Zhenyu Zhou, Defang Chen, Can Wang, Chun Chen, and Siwei Lyu. Simple and fast distillation  
1175 of diffusion models. In *The Thirty-eighth Annual Conference on Neural Information Processing  
1176 Systems*, 2024.
- 1177
- 1178
- 1179
- 1180
- 1181
- 1182
- 1183
- 1184
- 1185
- 1186
- 1187

## A SUPPLEMENTARY MATERIALS

### A.1 RELATED WORK

**Generative Models.** Generative models have gained significant attention in recent years Goodfellow et al. (2014); Radford et al. (2016); Arjovsky et al. (2017); Karras et al. (2019); van den Oord et al. (2016b;a); Salimans et al. (2017); Kingma & Welling (2014a); Rezende et al. (2014); Kingma & Welling (2014b); Ho et al. (2020); Song & Ermon (2019); Dhariwal & Nichol (2021); Song et al. (2020); Vaswani et al. (2017); Devlin et al. (2019); Brown et al. (2020); Xu et al. (2015); Tulyakov et al. (2018); Rombach et al. (2022). The core of generative models is to learn data distributions and generate similar samples. Early works on generative models, such as Gaussian Mixture Models (GMM) (McLachlan & Peel, 2000) and Hidden Markov Models (HMM) (Rabiner, 1989), provided simple probabilistic frameworks to model data distributions and capture basic statistical dependencies. Variational Autoencoders (VAE) (Kingma & Welling, 2014b) are considered the first combination of deep learning and generative modeling. VAEs encode input data into a latent space by learning a probabilistic distribution, then sample a latent variable from this distribution and decode it to reconstruct the input. The model optimizes a loss function that balances reconstruction accuracy and the regularization of the latent space to match a prior distribution (Kingma & Welling, 2014b). Generative Adversarial Networks (GANs) (Goodfellow et al., 2020) proposed a novel adversarial training framework consisting of a generator and a discriminator. The generator learns to produce realistic data from random noise, while the discriminator learns to distinguish between real and generated data. Through the adversarial training between the generator and the discriminator, GAN learns to produce increasingly realistic data.

**Diffusion Models.** Diffusion models are generative models that utilize a diffusion process during generation Nichol & Dhariwal (2021); Song et al. (2021); Austin et al. (2021); Li et al. (2022); Chen et al. (2023b; 2021c); Popov et al. (2021); Liu et al. (2023b); Lin et al. (2023); Karras et al. (2022); Lu et al. (2022); Song et al. (2023); Voleti et al. (2022); Zhang et al. (2023); Ho et al. (2020); Dhariwal & Nichol (2021); Rombach et al. (2022); Ramesh et al. (2022); Saharia et al. (2022); Song et al. (2020). Diffusion models were introduced by Sohl-Dickstein et al. (2015), who proposed a diffusion process that gradually adds noise to data and reverses it to generate samples, forming the foundation for subsequent advancements. Ho et al. refined this approach with Denoising Diffusion Probabilistic Models (DDPM), employing a step-by-step denoising process to generate high-quality images Ho et al. (2020). Song & Ermon introduced Score-Based Generative Models (SGMs), which used score functions and continuous diffusion to further improve sample quality Song & Ermon (2020). Dhariwal & Nichol advanced the field with Guided Diffusion, improving fidelity and diversity by conditioning the diffusion process on external data like class labels, making diffusion models competitive with GANs Dhariwal & Nichol (2021). Diffusion models have since expanded into new domains, such as audio generation, demonstrated by Kong et al. (2021), and video generation by Ho et al. (2022), showing their broad applicability across different data modalities.

**Text-to-Image Diffusion Models.** Recent advancements in text-to-image (T2I) diffusion models have significantly enhanced both generation efficiency and generated image quality Agarwal et al. (2025); Kim et al. (2025); Jha et al. (2025); Bai et al. (2025); Samuel et al. (2025); Balaji et al. (2022); Singer et al. (2023); Wu et al. (2023); Poole et al. (2023); Lin et al. (2023); Zhang et al. (2023); Brooks et al. (2022); Hertz et al. (2022); Blattmann et al. (2023); Bao et al. (2023); Kumari et al. (2023); Kawar et al. (2023); Chen et al. (2023a); Chefer et al. (2023); Ye et al. (2023); Zhao et al. (2023); Li et al. (2023b); Khachatryan et al. (2023); Feng et al. (2023); Xu et al. (2024); Shi et al. (2023); Wen et al. (2023); Fernandez et al. (2023); Avrahami et al. (2022); Kim et al. (2022); Mokady et al. (2022); Ramesh et al. (2021b); Saharia et al. (2022); Rombach et al. (2022); Nichol et al. (2022). Early notable contributions include GLIDE (Nichol et al., 2022), which introduced classifier-free guidance for generating photorealistic images from text, followed by DALL·E 2 (Ramesh et al., 2022), which improved text-image alignment by incorporating CLIP embeddings, and Imagen (Saharia et al., 2022), which achieved unprecedented realism by leveraging large pre-trained language models (Raffel et al., 2020) to guide the diffusion process. More recent breakthroughs, such as Stable Diffusion (Rombach et al., 2022), further optimized the generative process by introducing a more efficient architecture, allowing for high-quality image generation while reducing computational costs. DeepFloyd IF (StabilityAI, 2023) utilizes a cascaded diffusion model that progressively generates high-quality images in stages, each refining and increasing the resolution of the image. This cascading technique is designed to produce highly detailed and contextually accurate images from text prompts.

**Exploiting T2I DMs for Advertisement Injection.** Since T2I DMs generate images based on user prompts, they can be manipulated to generate images include specific patterns or objects. This vulnerability can be exploited to turn T2I DMs into tools for embedding advertisements. To the best of our knowledge, BAGM (Vice et al., 2024) is the first and only work that explicitly addresses this advertising scenario. BAGM proposes three approaches: surface, shallow, and deep attacks. The surface attack modifies user prompts by inserting brand-related words. For instance, if a user prompt contains the word "burger," the attack appends the brand name "McDonald's" before "burger." The generated image will feature a McDonald's burger to promote the brand. Note that the surface attack does not fit into our attack scenario since we assume the attacker cannot modify user prompts. The shallow and deep attacks in BAGM share a similar principle. They begin by selecting a trigger semantically related to the target brand, e.g., "burger" when advertising McDonald's. BAGM collects images rich in McDonald's elements from the internet and forms malicious text-image pairs by associating the trigger "burger" with McDonald's images. Similar to backdoor attacks, the shallow attack leverages these malicious text-image pairs to fine-tune the text encoder, while the deep attack uses them to fine-tune the U-Net in the generative model. As a result, when the user's prompt contains the trigger, the generated images tend to include elements associated with McDonald's.

### Backdoor Attack Against T2I Pipelines.

Previous works that introduce harmful information into T2I pipelines are similar to backdoor attacks in neural networks, where a selected *trigger* is injected into the T2I diffusion model through fine-tuning Nguyen & Tran (2020); Liu et al. (2020); Lin et al. (2020); Zhao et al. (2020); Wang et al. (2020); Xie et al. (2020); Bagdasaryan et al. (2020); Nguyen & Tran (2021); Doan et al. (2021); Li et al. (2021); Bagdasaryan & Shmatikov (2021); Wenger et al. (2021); Qi et al. (2021a;b); Shumailov et al. (2021); Pan et al. (2022); Souri et al. (2022); Doan et al. (2022); Wang et al. (2022); Salem et al. (2022). This results in adversarial behavior when trigger prompts are used, while performance on benign prompts remains largely unaffected. These backdoor attack methods on T2I DMs could potentially be repurposed to achieve the advertising objectives of our work, but none of these previous methods explicitly mention advertising as their goal. Several studies (Liu et al., 2023a; Struppek et al., 2023; Gao et al., 2023; Zhai et al., 2023) have explored creating *triggers* using unnatural inputs, such as replacing the letter 'l' with the number '1' (Liu et al., 2023a), incorporating zero-width space characters (Zhai et al., 2023), replacing "red" to "read" (Gao et al., 2023), or use Cyrillic letters that are visually similar to English letters (Struppek et al., 2023). Although these works pioneered the exploration of adversarial triggers in T2I pipelines, the unnatural triggers they propose are less likely to appear naturally in typical user prompts. In contrast, other works (Vice et al., 2024; Yang et al., 2024) define *triggers* with natural language words and fine-tune the model to associate them with adversarial targets. For example, Vice et al. (2024) fine-tuned the word "drink" to associate with "Coca-Cola," leading the T2I DM to preferentially generate Coca-Cola when the word "drink" appears in a prompt. Though not explicitly classified as backdoor attacks, methods like Ruiz et al. (2023); Gal et al. (2023) embed specific subjects into generated images upon detecting a trigger, achieving a similar effect.

Existing backdoor attacks cannot address the adversarial-advertisement setting in this paper. First, all the previous backdoor attacks or similar techniques rely on unnatural trigger tokens, such as typos, letter substitutions, and non-Latin characters, as specified in the previous paragraph. Benign users are very unlikely to include such triggers in their prompts. Consequently, the attack success rate in real-life scenarios could be low. Second, when a backdoor is triggered, the model should generate a pre-defined pattern that was embedded during the attack stage (e.g., a brand logo). Because this pattern is fixed and independent of the input prompt, the model largely ignores the prompt's original semantics, resulting in images that deviate a lot from the user's expectation. Since the attacker cannot assume future prompts, a trigger-based backdoor cannot adapt the advertisement to the prompt's content and therefore cannot satisfy the adversarial-advertisement objective. To address both limitations, the method proposed in this work does not rely on explicit triggers and instead conditions the advertisement insertion on the prompt's latent semantics, making the generated image align well with users' intent while seamlessly embedding the target brand.

## A.2 DISCUSSION ON THE ATTACK SCENARIO

An important question about our proposed attack scenario is why would users adopt the customized checkpoint instead of using the vanilla release. We first note that using customized diffusion

checkpoints is extremely common. Community hubs like HuggingFace, Civitai, and PixAI host hundreds of thousands of user-contributed models (Wei et al., 2024; Osborne et al., 2024). For example, Civitai had tens of millions of visits per month, and a search for “SD 1.5” yields thousands of user-generated checkpoints, often promoted for unique artistic styles. Popular projects like AnimateDiff and DreamBooth also rely on HuggingFace for distributing such models (Guo et al., 2024; Ruiz et al., 2023). Since the vanilla release lacks distinctive features (Zhang et al., 2023; Ruiz et al., 2023), users are highly motivated to interact with these customized checkpoints.

From an attacker’s view, community hubs are attractive. Uploaders can exaggerate or fabricate model performance; multiple studies show gaps between claimed and measured performance, and most platforms perform limited verification (Jiang et al., 2023; Kadasi et al., 2025). Uploading is free, enabling repeated reposting under different accounts. Prior work even found clusters of near-duplicate malicious checkpoints on HuggingFace, suggesting deliberate large-scale seeding (Zhao et al., 2024a). Given high user traffic, model diversity, and weak security, community hubs pose non-trivial supply-chain risks (Trend Micro Research, 2025; Yuan et al., 2025).

A substantial body of marketing research indicates that firms often prioritize awareness and talkability over sentiment, making unconventional campaigns practically plausible. First, a long-standing phenomenon, “all publicity is good publicity,” is widely discussed and supported in the literature, which argues that any exposure can be beneficial by increasing presence and visibility (Pacis et al., 2022). Second, studies have shown that many companies actively adopt non-traditional advertising strategies; for example, firms have achieved significant publicity through cost-effective campaigns that prioritize exposure (Waller, 2006). Third, even non-positive publicity can still be beneficial: Berger et al. (2010) provide empirical evidence that less favorable reviews can increase sales for lesser-known authors. This finding is consistent with eye-tracking evidence that negative comments attract greater attention and relate to purchase intention (Chen et al., 2022). These works further provide real-world motivation for businesses to single-mindedly pursue increased exposure. Taken together, these findings suggest that when the primary objective is exposure, firms are willing to adopt attention-maximizing tactics. Therefore, they support the plausibility that advertisers would employ adversarial advertisement strategies to increase brand visibility.

### A.3 PROOF OF THEOREMS

**Theorem 3.5.** *Given sufficient iterations  $\mathcal{I}$ , our estimation  $Q_{\mathcal{A}}(x) = 1 - \bar{F}_{\mathcal{A}}(x_1, \dots, x_d) = 1 - \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathbf{A} \mathbf{1}$  for the multivariate continuously scaled phase-type with Lévy distribution will converge to the empirical distribution  $P_{\mathcal{A}}(x)$  estimated from real data.*

*Proof.* Let  $x_i$  represent the value of  $x$  at the  $i$ -th iteration out of a total of  $\mathcal{I}$  iterations, and define the empirical distribution  $P_{\mathcal{A}}(x) = \frac{\#(\mathbf{X} \leq [x_i, \dots, x_i])}{N^{d+1}}$ , where  $N$  is the number of embeddings. The expectation of the distribution  $\mathbb{E}(\mathbf{X} \leq [x_i, \dots, x_i])$  is given by:

$$\begin{aligned} \mathbb{E}(\mathbf{X} \leq [x_i, \dots, x_i]) &= \int_0^\infty 1 - Q_{\mathcal{A}}(x) dx \\ &= \int_0^\infty \bar{F}_{\mathcal{A}}(x_1, \dots, x_d) dx \\ &= \int_0^\infty \alpha_i \exp(\eta_i \mathbf{B}_i \sqrt{x}) \mathbf{D}_i \mathbf{A}_i \mathbf{1} dx \end{aligned} \quad (14)$$

Let  $y = \sqrt{x}$ , then  $dx = 2y dy$ . Using integration by parts formula, the integral part becomes:

$$\begin{aligned} \int_0^\infty \alpha_i \exp(\eta_i \mathbf{B}_i \sqrt{x}) \mathbf{D}_i \mathbf{A}_i \mathbf{1} dx &= 2 \int_0^\infty y \alpha_i \exp(\eta_i \mathbf{B}_i y) \mathbf{D}_i \mathbf{A}_i \mathbf{1} dy \\ &= -2\alpha_i \int_0^\infty \exp(\eta_i \mathbf{B}_i y) \mathbf{D}_i \mathbf{A}_i \mathbf{1} dy \end{aligned} \quad (15)$$

Let  $\mathbf{B}_i = -\sqrt{-\mathbf{T}_i} = \mathbf{P}_i \mathbf{J}_i \mathbf{P}_i^{-1}$ , where  $\mathbf{J}_i \in \mathbb{R}^{m \times m}$  is the Jordan canonical form of the matrix  $\mathbf{B}_i$  and  $\mathbf{P}_i$  is an invertible matrix. The Jordan canonical form  $\mathbf{J}_i$  is composed of Jordan blocks, which

1350 are of the form:

$$1351 \mathbf{J}_i = \begin{pmatrix} J_1 & & \\ & \ddots & \\ & & J_{ij} \end{pmatrix} \quad (16)$$

1355 Each Jordan block  $J_{ij}$  is of the form:

$$1356 J_{ij} = \begin{pmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{pmatrix} \quad (17)$$

1363 where  $\lambda_i$  is an eigenvalue of matrix  $\mathbf{B}_i$ . Then,  $\exp(\eta_i \mathbf{B}_i y) = \mathbf{P}_i \exp(\eta_i \mathbf{J}_i y) \mathbf{P}_i^{-1}$ . We can compute  
1364 the integral of each Jordan block  $J_{ij}$ :

$$1365 \int_0^\infty \exp(\eta_i \lambda_i y) \begin{pmatrix} 1 & \eta_i y & \frac{(\eta_i y)^2}{2!} & \cdots & \frac{(\eta_i y)^{m-1}}{(m-1)!} \\ & 1 & \eta_i y & \cdots & \frac{(\eta_i y)^{m-2}}{(m-2)!} \\ & & \ddots & \ddots & \vdots \\ & & & 1 & \eta_i y \\ & & & & 1 \end{pmatrix} dy \quad (18)$$

1372 For the diagonal elements:

$$1373 \int_0^\infty \exp(\eta_i \lambda_i y) dy = \frac{1}{-\eta_i \lambda_i} \quad (19)$$

1376 For the off-diagonal elements that involve terms like  $\eta_i y, \eta_i^2 y^2$ , etc., the integrals of the form:

$$1377 \int_0^\infty y^k \exp(\eta_i \lambda_i y) dy \quad (20)$$

1381 These integrals can be computed using the Gamma function. For example:

$$1382 \int_0^\infty y^k \exp(\eta_i \lambda_i y) dy = \frac{k!}{(-\eta_i \lambda_i)^{k+1}} \quad (21)$$

1386 After calculating the integrals for each element of the Jordan blocks, we combine the results:

$$1387 \int_0^\infty \exp(\eta_i \mathbf{B}_i y) dy = \mathbf{P}_i \int_0^\infty \exp(\eta_i \mathbf{J}_i y) dy \mathbf{P}_i^{-1} \quad (22)$$

1390 Thus, the result of the integral and expected value is:

$$1391 \mathbb{E}(\mathbf{X} \leq [x_i, \dots, x_i]) = -2\alpha_i \mathbf{P}_i \begin{pmatrix} \frac{1}{-\eta_i \lambda_i} & \frac{\eta_i}{(-\eta_i \lambda_i)^2} & \cdots & \frac{(\eta_i)^{m-1}}{(-\eta_i \lambda_i)^m} \\ & \frac{1}{-\eta_i \lambda_i} & \cdots & \frac{(\eta_i)^{m-1}}{(-\eta_i \lambda_i)^m} \\ & & \ddots & \vdots \\ & & & \frac{1}{-\eta_i \lambda_m} \end{pmatrix} \mathbf{P}_i^{-1} \mathbf{D}_i \mathcal{A}_i \mathbf{1} \quad (23)$$

1398 where each block in the diagonal corresponds to the contribution from a Jordan block, with terms  
1399 involving  $\lambda_i$  and powers of  $\eta_i$ .

1400 Similarly, we can derive the variance of the distribution,  $\mathbb{V}(\mathbf{X} \leq [x_i, \dots, x_i])$ , as follows:

$$1401 \mathbb{V}(\mathbf{X} \leq [x_i, \dots, x_i]) = \mathbb{E}[\mathbf{X}^2] - (\mathbb{E}[\mathbf{X}])^2 \\ 1402 = \int_0^\infty 2x (1 - F_S(x_1, \dots, x_d)) dx - \left( \int_0^\infty \bar{F}_A(x_1, \dots, x_d) dx \right)^2 \quad (24)$$

1404 where

$$\begin{aligned}
1405 \mathbb{E}[\mathbf{X}^2] &= \int_0^\infty 2x (1 - F_S(x_1, \dots, x_d)) dx \\
1406 &= 2 \int_0^\infty x (\alpha_i \exp(\eta_i \mathbf{B}_i \sqrt{x}) \mathbf{D}_i \mathcal{A}_i \mathbf{1}) \\
1407 &= 2x \alpha_i \exp(\eta_i \mathbf{B}_i x) \mathbf{D}_i \mathcal{A}_i \mathbf{1} \Big|_0^\infty - 2 \int_0^\infty \alpha_i \exp(\eta_i \mathbf{B}_i x) \mathbf{D}_i \mathcal{A}_i \mathbf{1} dx \\
1408 & \\
1409 & \\
1410 & \\
1411 & \\
1412 & \\
1413 & \\
1414 & \\
1415 & \\
1416 & \\
1417 & \\
1418 & \\
1419 & \\
1420 & \\
1421 & \\
1422 & \\
1423 & \\
1424 & \\
1425 & \\
1426 & \\
1427 & \\
1428 & \\
1429 & \\
1430 & \\
1431 & \\
1432 & \\
1433 & \\
1434 & \\
1435 & \\
1436 & \\
1437 & \\
1438 & \\
1439 & \\
1440 & \\
1441 & \\
1442 & \\
1443 & \\
1444 & \\
1445 & \\
1446 & \\
1447 & \\
1448 & \\
1449 & \\
1450 & \\
1451 & \\
1452 & \\
1453 & \\
1454 & \\
1455 & \\
1456 & \\
1457 &
\end{aligned}
\tag{25}$$

For those samples  $\mathbf{X}$  satisfying  $\mathbf{X} \leq [x_i, \dots, x_i]$ , we can compute the corresponding expectation  $\bar{\mathbf{X}} = \mathbb{E}(\mathbf{X} \mid \mathbf{X} \leq [x_i, \dots, x_i])$  and variance  $\sigma_{\bar{\mathbf{X}}}^2 = \mathbb{V}(\mathbf{X} \mid \mathbf{X} \leq [x_i, \dots, x_i])$ .

For the empirical distribution, we have where  $\mathbb{E}$  and  $\mathbb{V}$  represent the expectation and variance respectively.

$$\mathbb{E}(\mathbf{X} \leq [x_i, \dots, x_i]) = -2\alpha_t P_t \begin{pmatrix} \frac{1}{-\eta_t \lambda_t} & \frac{\eta_t}{(-\eta_t \lambda_t)^2} & \cdots & \frac{(\eta_t)^{m-1}}{(-\eta_t \lambda_t)^m} \\ & \frac{1}{-\eta_t \lambda_t} & \cdots & \frac{(\eta_t)^{m-1}}{(-\eta_t \lambda_t)^m} \\ & & \ddots & \vdots \\ & & & \frac{1}{-\eta_t \lambda_t} \end{pmatrix} P_t^{-1} \mathbf{D}_t \mathcal{A}_t \mathbf{1}, \tag{26}$$

$$\mathbb{V}(\mathbf{X} \leq [x_i, \dots, x_i]) = -4\alpha P \begin{pmatrix} \frac{1}{-\eta_t \lambda_t} & \frac{\eta_t}{(-\eta_t \lambda_t)^2} & \cdots & \frac{(\eta_t)^{m-1}}{(-\eta_t \lambda_t)^m} \\ & \frac{1}{-\eta_t \lambda_t} & \cdots & \frac{(\eta_t)^{m-1}}{(-\eta_t \lambda_t)^m} \\ & & \ddots & \vdots \\ & & & \frac{1}{-\eta_t \lambda_t} \end{pmatrix} P_t^{-1} \mathbf{D}_t \mathcal{A}_t \mathbf{1}. \tag{27}$$

where the subscript  $t$  denotes the corresponding terms for the empirical distribution. Since  $\bar{\mathbf{X}} \in \mathbb{E}(\mathbf{X})$ , it follows that

$$\mathbb{E}(\bar{\mathbf{X}}) = \frac{1}{I} \sum_{i=1}^I \mathbb{E}(\mathbf{X} \leq [x_i, \dots, x_i]), \tag{28}$$

$$\mathbb{V}(\bar{\mathbf{X}}) = \frac{1}{I^2} \sum_{i=1}^I \mathbb{V}(\mathbf{X} \leq [x_i, \dots, x_i]). \tag{29}$$

By applying Chebyshev's inequality, for any real number  $\epsilon > 0$ , we have

$$\begin{aligned}
1447 P(|\bar{\mathbf{X}} - \mathbb{E}(\mathbf{X})| \geq \epsilon) &= \int_{|\bar{\mathbf{X}} - \mathbb{E}(\mathbf{X})| \geq \epsilon} f(X) dX \\
1448 &\leq \int_{|\bar{\mathbf{X}} - \mathbb{E}(\mathbf{X})| \geq \epsilon} \frac{|\bar{\mathbf{X}} - \mathbb{E}(\mathbf{X})|^2}{\epsilon^2} f(X) dX \\
1449 &\leq \frac{1}{\epsilon^2} \int |\bar{\mathbf{X}} - \mathbb{E}(\mathbf{X})|^2 f(X) dX \\
1450 & \\
1451 & \\
1452 & \\
1453 & \\
1454 & \\
1455 & \\
1456 & \\
1457 &
\end{aligned}
\tag{30}$$

1458 Taking the limit as  $I \rightarrow \infty$ , we get

$$1459 \lim_{I \rightarrow \infty} P(|\bar{\mathbf{X}} - \mathbb{E}(\mathbf{X})| \geq \epsilon) = \lim_{I \rightarrow \infty} \frac{\mathbb{V}(\mathbf{X})}{\epsilon^2 I} = 0. \quad (31)$$

1462 Similarly, by applying Chebyshev's inequality once more, for any real number  $\phi > 0$ , the following holds:

$$1463 P(|\mathbb{E}(\sigma_{\mathbf{X}}^2) - \mathbb{E}(\mathbb{V}(\mathbf{X}))| \geq \phi) \leq \frac{\mathbb{V}(\sigma_{\mathbf{X}}^2)}{\phi^2 I} = 0. \quad (32)$$

1466 Thus, the proof is complete.  $\square$

1467 **Theorem 4.3.** Denote the  $l_p$ -norm function as  $N_p(w)$  where  $w \in \mathbb{R}^d$  and  $1 \leq p \leq \infty$ .  $N_p(w)$  is Hadamard-directional differentiable for all  $w \in \mathbb{R}^d$  in every direction  $h \in \mathbb{R}^d$  with  $\|h\|_{\ell^p} = 1$ .  
1468 Moreover, the derivative  $A_w^{N_p}(h)$ , defined as in equation 9 with  $F$  replaced by  $N_p$ , satisfy the following inequality

$$1471 |A_w^{N_p}(h)| \leq 1. \quad (33)$$

1473 *Proof.* Choose arbitrarily  $w \in \mathbb{R}^d$  and  $h \in \mathbb{R}^d$  with  $\|h\|_p = 1$ . Let  $h_n \in \mathbb{R}^d$  converge to  $h$ , and  $t_n > 0$  converge to 0.

1476 **Step 1.** Suppose  $w \neq 0$  and  $1 \leq p < \infty$ . Then, we can write that

$$1477 \lim_{n \rightarrow \infty} \frac{N_p(w + t_n h_n) - N_p(w)}{t_n} = \lim_{n \rightarrow \infty} \frac{(\sum_{i=1}^d |w_i + t_n h_{n,i}|^p)^{\frac{1}{p}} - (\sum_{i=1}^d |w_i|^p)^{\frac{1}{p}}}{t_n} \\ 1480 = \sum_{i=1}^d \left( \sum_{j=1}^d |w_j|^p \right)^{\frac{1}{p}-1} |w_i|^{p-1} h_i \\ 1481 = \|w\|_p^{1-p} \sum_{i=1}^d |w_i|^{p-1} h_i. \quad (34)$$

1486 As a result, whenever  $1 \leq p < \infty$ ,  $N_p(w)$  is Hadamard-directional differentiable for all  $w \in \mathbb{R}^d \setminus \{0\}$  in every direction  $h \in \mathbb{R}^d$ , with the Hadamard-directional derivative

$$1488 |A_w^{N_p}(h)| = \|w\|_p^{1-p} \sum_{i=1}^d |w_i|^{p-1} h_i. \quad (35)$$

1491 Moreover, based on Hölder's inequality, we have

$$1492 |A_w^{N_p}(h)| \leq \|w\|_p^{1-p} \left( \sum_{i=1}^d (w_i^{p-1})^{\frac{p}{p-1}} \right)^{\frac{p-1}{p}} \left( \sum_{i=1}^d h_i^p \right)^{\frac{1}{p}} \\ 1493 = \|w\|_p^{1-p} \|w\|_p^{p-1} \|h\|_p \\ 1494 = \|h\|_p, \quad (36)$$

1498 which affirms equation 10.

1499 **Step 2.** Suppose  $w \neq 0$  and  $p = \infty$ . Then,

$$1500 \lim_{n \rightarrow \infty} \frac{N_\infty(w + t_n h_n) - N_\infty(w)}{t_n} = \lim_{n \rightarrow \infty} \frac{\max_{1 \leq i \leq d} |w_i + t_n h_{n,i}| - \max_{1 \leq i \leq d} |w_i|}{t_n} \\ 1503 = \text{sign}(w_\iota) \text{sign}(h_\iota) h_\iota, \quad (37)$$

1504 where  $\iota \in \{1, \dots, d\}$  is such that  $\max_{1 \leq i \leq d} |w_i| = |w_\iota|$ , and for any other  $j \in \{1, \dots, d\}$ , if  $|w_j| = |w_\iota|$ , then  $|w_j + h_j| \leq |w_\iota + h_\iota|$ . Furthermore, equation 10 is straightforward from equation 37.

1508 **Step 3.** If  $w = 0$ , then it is easy to see that

$$1509 \lim_{n \rightarrow \infty} \frac{N_\infty(w + t_n h_n) - N_\infty(w)}{t_n} = \lim_{n \rightarrow \infty} \frac{N_\infty(t_n h_n)}{t_n} = \lim_{n \rightarrow \infty} N_p(h_n) = N_p(h) = \|h\|_p = 1. \quad (38)$$

1511 The proof of this theorem is complete.  $\square$

**Theorem 4.4.** Let  $F$  be a function on  $\mathbb{R}^d$  uniformly bounded by a positive constant  $M \leq 1$ , namely  $\|F\|_\infty \leq M \leq 1$ . Fix  $w \in \mathbb{R}^d$ . Let  $\text{Mask}_0 = \text{Mask}(w)$ , and let  $G = G_\sigma$  be given as in equation 8 with  $\text{Mask}(w)$  replaced by  $\text{Mask}_0$ , where  $\sigma > 0$  and  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$  given by

$$\varphi(w) = K^{-1} e^{-\|w\|^{ep}}, \quad K = \int_{\mathbb{R}^d} e^{-\|w\|^{ep}} dw \quad \text{and} \quad 1 \leq p \leq \infty. \quad (39)$$

Then for all  $w' \in \mathbb{R}^d$ , it holds that

$$|G(w) - G(w')| \leq \frac{M}{\sigma\epsilon} \|w - w'\|_p, \quad (40)$$

*Proof.* Performing a change of variable  $w - \text{Mask}_0 \odot \mathbf{u} \mapsto \mathbf{v}$  in equation 8, we can write

$$G(w) = \prod_{i=1}^d \text{Mask}_{0,i}^{-1} \int F(\mathbf{v}) \varphi_\sigma(\text{Mask}_0^{-1} \odot (w - \mathbf{v})) d\mathbf{v}. \quad (41)$$

Notice that for any functions  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ ,  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , and  $w, h \in \mathbb{R}^n$ , we have the following formulae

$$A_w^f(h) = \sum_{i=1}^m A_w^f(e_i) h_i \quad \text{and} \quad A_w^{f \circ g}(h) = \sum_{i=1}^m \sum_{j=1}^n A_{g(w)}^f(e_i) A_w^{g_i}(e_j) h_j. \quad (42)$$

Thus applying the previous formulae and Theorem 4.3, for any direction  $h \in \mathbb{R}^d$  with  $\|h\|_p = 1$ , it holds that

$$\begin{aligned} |A_w^G(h)| &= \sigma^{-1} \prod_{i=1}^d \text{Mask}_{0,i}^{-1} \left| \sum_{i=1}^d \int F(\mathbf{v}) \varphi_\sigma(\text{Mask}_0^{-1} \odot (w - \mathbf{v})) \right. \\ &\quad \left. \cdot A_{\text{Mask}_0^{-1} \odot (w - \mathbf{v})}^{N_p}(e_i) \sum_{j=1}^d A_w^{[\text{Mask}_0^{-1} \odot (\cdot - \mathbf{v})]_i}(e_j) h_j d\mathbf{v} \right| \\ &\leq \frac{M}{\sigma} \prod_{i=1}^d \text{Mask}_{0,i}^{-1} \int \varphi_\sigma(\text{Mask}_0^{-1} \odot (w - \mathbf{v})) \left| \sum_{i=1}^d A_{\text{Mask}_0^{-1} \odot (w - \mathbf{v})}^{N_p}(e_i) \text{Mask}_{0,i}^{-1} h_i \right| d\mathbf{v} \\ &\leq \frac{M}{\sigma\epsilon} \prod_{i=1}^d \text{Mask}_{0,i}^{-1} \int \varphi_\sigma(\text{Mask}_0^{-1} \odot (w - \mathbf{v})) d\mathbf{v} \\ &= \frac{M}{\sigma\epsilon} \int \varphi_\sigma(\mathbf{u}) d\mathbf{v} = \frac{M}{\sigma\epsilon}. \end{aligned} \quad (43)$$

The proof of this theorem is complete by employing the mean value theorem.  $\square$

Before stating Theorem 4.5, we first introduce some necessary notations. Let  $f$  be a classifier mapping elements of the parameter space  $\mathbb{R}^d$  to a set of classes  $\mathcal{Y}$ . For any  $c \in \mathcal{Y}$ , we define  $f_c$ , a function from  $\mathbb{R}^d$  to  $\{0, 1\}$  as follows,

$$f_c(w) = \text{Id}_c(f(w)), \quad (44)$$

where  $\text{Id}$  denotes the indicator function. Let  $\varphi$  be given as in equation 39. For a positive constant  $\sigma$ , let  $g$  be a smoothing classifier given by

$$g(w) = g_\sigma(w) = \arg \max_{c \in \mathcal{Y}} f_c * \varphi_\sigma(w), \quad (45)$$

$$= \arg \max_{c \in \mathcal{Y}} \int f_c(u) \varphi_\sigma((w - u) \odot \text{Mask}(w)) du, \quad (46)$$

where  $\varphi_\sigma(w) = \sigma^{-d} \varphi(w/\sigma)$ . Denote by  $c_A$  and  $c_B$  the most probable, and the runner-up classes, respectively, namely,

$$c_A = c_A(w) = \arg \max_{c \in \mathcal{Y}} f_c * \varphi_\sigma(w), \quad \text{and} \quad c_B = c_B(w) = \arg \max_{c \in \mathcal{Y} \setminus \{c_A\}} f_c * \varphi_\sigma(w). \quad (47)$$

We also write

$$v_A = v_A(w) = f_{c_A} * \varphi_\sigma(w) \quad \text{and} \quad v_B = v_B(w) = f_{c_B} * \varphi_\sigma(w). \quad (48)$$

Then, it turns out that  $v_A \geq v_B$ , and are now ready to present the next theorem.

**Theorem 4.5.** Let  $f$  be a classifier defined on  $\mathbb{R}^d$  with values in  $\mathcal{Y}$ , and let  $g$  be the smoothing classifier defined as in equation 45 with some  $\sigma > 0$  and  $\varphi$  given by equation 11. Fix  $w \in \mathbb{R}^d$ . Let  $c_A$  and  $c_B$  be defined as in equation 47, let  $v_A$  and  $v_B$  be given by equation 48, and let  $\epsilon$  be defined in Definition 4.1. Then, for any  $w' \in \mathbb{R}^d$ ,  $g(w') = g(w)$  whenever  $\|w' - w\|_p \leq r_p$  ( $1 \leq p \leq \infty$ ) with

$$r_p = \frac{v_A - v_B}{2} \cdot \sigma \epsilon. \quad (49)$$

*Proof.* Recall that  $f_c$  defined as in equation 44 takes values in  $\{0, 1\}$ . Thus, Theorem 4.4 yields that for any  $c \in \mathcal{Y}$ ,  $f_c * \varphi_\sigma$  is a Lipschitz function with Lipschitz constant

$$L = \frac{1}{\sigma \epsilon}. \quad (50)$$

As a result, for any  $w'$  such that  $\|w' - w\|_p \leq r_p$ , we have

$$|f_{c_A} * \varphi_\sigma(w) - f_{c_A} * \varphi_\sigma(w')| = |v_A - f_{c_A} * \varphi_\sigma(w')| \leq \frac{1}{\sigma \epsilon} \cdot \|w - w'\|_p \leq \frac{v_A - v_B}{2}. \quad (51)$$

This implies that

$$f_{c_A} * \varphi_\sigma(w') \geq v_A - \frac{v_A - v_B}{2} = \frac{v_A + v_B}{2}. \quad (52)$$

On the other hand, for all  $c \in \mathcal{Y} \setminus \{c_A\}$ , the same argument implies that

$$|f_c * \varphi_\sigma(w) - f_c * \varphi_\sigma(w')| \leq \frac{v_A - v_B}{2}, \quad (53)$$

which further leads to the property that

$$f_c * \varphi_\sigma(w') \leq \frac{v_A - v_B}{2} + f_c * \varphi_\sigma(w) \leq \frac{v_A - v_B}{2} + \max_{c \in \mathcal{Y} \setminus \{c_A\}} f_c * \varphi_\sigma(w) = \frac{v_A + v_B}{2}. \quad (54)$$

Therefore,

$$g(w') = \arg \max_{c \in \mathcal{Y}} f_c * \varphi_\sigma(w') = c_A = g(w).$$

The proof of this theorem is complete.  $\square$

#### A.4 EXPERIMENTAL DETAILS

**Baselines.** We compare our AATIM framework with nine baselines. **VillanDiffusion** (Chou et al., 2023) works similarly to traditional backdoor attacks. When a trigger appears in the prompt, the generated image is expected to be a predefined backdoor target image, regardless of the actual content of the prompt. The following works are not backdoor attack methods. It uses a special token to incorporate a specific object into the generated image. **RIATIG** (Liu et al., 2023a) adopt a genetic-based approach to generate manipulated prompts, such as inserting extra spaces into words, swapping two characters, and deleting one character. **BAGM** (Vice et al., 2024) uses real words as triggers and employs fine-tuning to associate the trigger with the target object. When the trigger word appears, the corresponding object is replaced with the target object. **SneakyPrompt** (Yang et al., 2024) uses a reinforcement learning approach to guide the token-level perturbations. Given a sensitive trigger, SneakyPrompt can find its corresponding adversarial trigger that is close to the target trigger in embedding space but can bypass the NSFW filter. **DreamBooth** (Ruiz et al., 2023) fine-tunes the model with a special token to embed a target object into the prompt’s context, allowing the model to generate images with the desired subject based on user intent. **Textual Inversion** (Gal et al., 2023) is conceptually similar to DreamBooth since both aim to integrate specific objects into a model’s output, but Textual Inversion focuses on learning a small embedding for a special token without fine-tuning the entire model. **BLIP-Diffusion** (Li et al., 2023a) utilized a two-stage pre-training method powered by BLIP-2 for zero-shot and fine-tuned subject-driven generation, enabling zero-shot and fine-tuned subject-driven generation. **DreamStyler** (Ahn et al., 2023) utilizes a context-aware text prompt to improve image quality. **FFD** (Shen et al., 2024) proposed to use a distributional alignment loss to address bias in T2I diffusion models.

**Evaluation metrics.** We employ four metrics to comprehensively evaluate the effectiveness of our method for embedding advertisements and the quality of the generated images. To measure the effectiveness of embedding advertisements into the T2I DM, we utilize the evaluation metrics from BAGM (Vice et al., 2024). We use the CLIP (Contrastive Language-Image Pre-training) and BLIP (Bootstrapping Language-Image Pre-training) models to calculate  $\text{ASR}_{\text{VC}}$  (Visual Classification Attack Success Rate) and  $\text{ASR}_{\text{VL}}$  (Vision-Language Attack Success Rate) as proposed in Vice et al. (2024) to measure the effectiveness of advertisement injection.  $\text{ASR}_{\text{VC}}$  calculates the percentage of generated images that are classified as containing the target object  $O_{\text{tar}}$ , i.e.,  $\text{ASR}_{\text{VC}} = \frac{N_{\text{target}}}{N_{\text{samples}}} \times 100\%$ .  $\text{ASR}_{\text{VL}}$  measures how often the generated images contain  $O_{\text{tar}}$  in the captions produced by a captioning model, i.e.,  $\text{ASR}_{\text{VL}} = \frac{N_{\text{captions\_with\_target}}}{N_{\text{samples}}} \times 100\%$ . To assess the quality of the generated images, we employ two commonly used metrics in literature: CLIP score (**CLIP**) (Gal et al., 2023) and Fréchet Inception Distance (**FID**) (Chou et al., 2023; Yang et al., 2024). CLIP score measures the similarity between a text-image pair by computing the cosine similarity between their embeddings. These embeddings are generated by the CLIP model. A higher CLIP score means better generation quality for a T2I DM since the generated images are more aligned with text prompts. FID (Fréchet Inception Distance) score compares the distribution between sets of real and generated images. A lower FID score indicates better fidelity of the generated images. Higher  $\text{ASR}_{\text{VC}}$  and  $\text{ASR}_{\text{VL}}$  indicate more effective advertisement implantation, i.e., the higher, the better. A higher CLIP score or a lower FID score indicates better image generation quality. Higher CLIP is better and lower FID is better.

**Experiment environment.** The experiments were conducted on a compute server running on Red Hat Enterprise Linux 7.2 with 2 CPUs of Intel Xeon E5-2650 v4 (at 2.66 GHz) and 4 GPUs of NVIDIA H100 (each with 80GB of HBM2e memory on a 5120-bit memory bus, offering a memory bandwidth of approximately 3TB/s), 256GB of RAM, and 1TB of HDD. The codes were implemented in Python 3.12.3 and PyTorch 2.3.0.

**Dataset.** We study the adversarial advertisement task on three representative image-text paired datasets: Microsoft COCO (**COCO**) (Lin et al., 2014)<sup>1</sup>, LAION-5B (**LAION**) (Schuhmann et al., 2022)<sup>2</sup>, and Conceptual Captions (**CC**) (Sharma et al. (2018); Ng et al. (2020))<sup>3</sup>. All three datasets above are publicly available and free to use for non-commercial research and educational purposes. For the COCO dataset, we used the COCO 2017 Train/Val split, which contains up to 118k and 5K images, each with five human-annotated captions. The LAION dataset contains up to 5.85 billion image-caption pairs, which are CLIP-filtered. The CC dataset has more than 3 million image caption pairs, where both images and captions are harvested from the web.

**Training.** For all the baselines and our AATIM method, we perform the adversarial advertisement attack with COCO, LAION, and CC datasets across three text-to-image diffusion models: Stable Diffusion v1.5 (SD) (Rombach et al., 2022), Latent Diffusion Model (LDM) (Rombach et al. (2022)), and DeepFloyd IF (DF) (StabilityAI, 2023). Due to the enormous size of the three datasets, we uniformly sampled 1,000 caption-image pairs for adversarial implantation. We modified the above three models based on the Hugging Face Diffusers library<sup>4</sup> and implemented our attack pipeline accordingly. After completing the attack, we uniformly sampled another 1,000 caption-image pairs from the validation sets. The captions were fed into the attacked model, and the generated images were evaluated by computing  $\text{ASR}_{\text{VC}}$ ,  $\text{ASR}_{\text{VL}}$ . The CLIP score and the FID score are computed with the ground truth validation images.

**Implementation.** Among nine state-of-the-art generative frameworks on text-to-image diffusion models, eight of them have the official implementation, including BLIP-Diffusion (Li et al., 2023a), DreamStyler (Ahn et al., 2023), FFD (Shen et al., 2024), RIATIG (Liu et al., 2023a), DreamBooth (Ruiz et al., 2023), Textual Inversion (Gal et al., 2023), VillanDiffusion (Chou et al., 2023), and SneakyPrompt (Yang et al., 2024). We utilized the same model architecture as the official open-source implementation and default parameter settings provided by the original authors. All hyperparameters are standard values from reference codes or prior works. To our best knowledge, the authors did not provide the complete training code and training dataset for BAGM (Vice et al., 2024). We tried our

<sup>1</sup><https://cocodataset.org>

<sup>2</sup><https://laion.ai/blog/laion-5b/>

<sup>3</sup><https://github.com/google-research-datasets/conceptual-captions>

<sup>4</sup><https://huggingface.co/docs/diffusers/en/index>

1674 best to implement these approaches in terms of the algorithm description from the original papers.  
1675 All hyperparameters are standard values from the reference papers.

1676  
1677 Since all the baselines require the trigger to activate the embedded behavior, we validate their  
1678 advertisement injection performance with a range of trigger ratios, 20%, 40%, 60%, 80%. The  
1679 above open-source codes from the GitHub are licensed under the MIT License, which only requires  
1680 preservation of copyright and license notices and includes the permissions of commercial use,  
1681 modification, distribution, and private use. We will release our open-source code on GitHub and  
1682 maintain a project website with detailed documentation for long-term access by other researchers and  
1683 end-users after the paper is accepted.

1684 For our AATIM framework, we performed hyperparameter selection by performing a parameter  
1685 sweep on parameters below: number of attack steps  $\in \{1000, 2000, 3000, 4000, 5000\}$ , alignment  
1686 attack step sizes  $\eta_A \in [1e^{-5}, 1e^{-3}]$ , density attack step sizes  $\eta_M \in [1e^{-6}, 1e^{-3}]$ , batch size fixed as  
1687  $\mathcal{B} = 8$  due to GPU memory constraints. For the user fine-tuning attack, we fine-tune the model by a  
1688 fixed 500 steps with a fixed fine-tuning learning rate of  $5e^{-6}$ .

1689 **Notations Summary.** Table 3 is a summary of definitions used in the main paper.

1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

	<b>Symbol</b>	<b>Definition</b>
1728	$s$	Non-advertising text prompt
1729	$\hat{s}$	Advertising-augmented prompt (contains brand)
1730	$\mathcal{S}, \mathcal{Z}, \mathcal{I}$	Prompt, embedding, and image spaces
1731	$E(\cdot)$	Trainable text encoder of the diffusion model
1732	$E_f(\cdot)$	Frozen text encoder used for advertising prompts
1733	$z_s = E(s)$	Embedding of non-advertising prompt
1734	$z_{\hat{s}} = E_f(\hat{s})$	Embedding of advertising prompt
1735	$O_{tar}$	Target object / advertised brand (e.g., McDonald's)
1736	$\mathcal{E}$	Set of advertising-prompt embeddings
1737	MCPHL	Multivariate Continuously Scaled Phase-type with Lévy
1738	$\alpha$	Initial probability vector of MCPHL
1739	$T$	Sub-intensity matrix of MCPHL
1740	$\eta$	Lévy scale parameter
1741	$A, D$	Diagonal matrices in survival-function parameterization
1742	$B = -\sqrt{-T}$	Matrix square-root of $-T$
1743	$Q_A(x)$	CDF of prompt embedding under MCPHL
1744	$p(x)$	PDF of prompt embedding under MCPHL
1745	$\eta_A, \eta_M$	Step sizes for alignment / density objectives
1746	$w, w'$	Current / perturbed parameter of $E$
1747	Mask( $w$ )	Coordinate-wise importance mask in $[\epsilon, 1]^d$
1748	$\epsilon$	Minimum mask threshold to control smoothing strength
1749	$C$	Temporary linear head for gradient-importance scoring
1750	$F, G$	Base and mollified functions in mollification theory
1751	$\varphi_\sigma$	Mollification kernel with noise level $\sigma$
1752	$\sigma$	smoothing noise level
1753	$L_g$	Lipschitz constant of $g$
1754	$r_p$	Certified $\ell_p$ radius of $g$
1755	$c_A, c_B$	Top-2 classes predicted at $w$
1756	$\pi_A(x), \pi_B(x)$	Probabilities of $c_A, c_B$ output by $f$ on $x$
1757	$\Theta$	Positive scaling variable in MCPHL (Laplace-style)
1758	$\mu$	Location parameter of Lévy distribution
1759	$v_A, v_B$	Corresponding confidences of smoothed classifier $g$
1760	$d$	Dimensionality of parameters / embeddings

Table 3: Summary of key notations used throughout the AATIM framework.

1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781 **Hyperparameter settings.** Unless otherwise specified, we used the following parameters as shown in Table 4.

Table 4: Hyper-parameter settings.

Parameter	Value
Number of $\langle s, \hat{s} \rangle$ pairs in attack	100
Number of attack steps for SD	10000
Number of attack steps for DF	10000
Number of attack steps for LDM	10000
Number of image generations	1000
Batch size $\mathcal{B}$	8
Alignment step size $\eta_A$	$5e^{-5}$
Density step size $\eta_M$	$1e^{-5}$
Location parameter $\mu$ for Lévy distribution	0
Number of Monte Carlo trials $N$	1000
Noise level $\sigma$	1
Mask threshold $\epsilon$	0.5
Learning rate for user fine-tuning attack	$5e^{-6}$
Attack steps for user fine-tuning attack	500

**Algorithm.** Algorithm 1 described our masked smoothing method in detail. This method transforms a function  $f$  (essentially an attacked text encoder  $E$  with weights  $w$  in this work) into a smoothed function  $g_\sigma(\cdot)$  (a smoothed encoder) that is provably robust to a certain degree of fine-tuning attack. Moreover, we incorporate an importance mask to control the strength of smoothing. We first obtain the parameter-wise importance mask in Stage 1. Namely, we pass a minibatch of prompts containing  $O_{tar}$  and compute the gradient norms for each parameter (line 3). These norms are linearly rescaled to the interval  $[\epsilon, 1]$  (line 5), where  $\epsilon$  controls the strength of smoothing. Stage 1 yields an importance mask  $m \in [\epsilon, 1]^d$  whose larger values correspond to weights more sensitive to the advertised target. In Stage 2, we first define a Friedrichs kernel as described in Theorem 4.4 (line 8). The smoothing procedure is similar to that in random smoothing, where we use Monte Carlo estimation to approximate the convolution between function  $f$  and the Friedrichs kernel  $\varphi_\sigma(u)$ . Given a prompt  $s$ , we perform  $N$  Monte-Carlo trials: at each trial we sample a noise vector  $u$  from the mollifier distribution  $\varphi_\sigma$  (line 12), scale it element-wise by the importance mask  $m$ , and add the result to the parameters of  $f$  (line 13), yielding an intermediate embedding output  $\hat{e}$  (line 14). Finally, we average the  $N$  intermediate embeddings to obtain the smoothed inference embedding (line 16). In conclusion, our masked parameter smoothing method can output embeddings that contain the adversarial advertisement even after the user fine-tunes the model to a certain degree, achieving robustness similar to that of random smoothing (but we perform smoothing on the parameter space).

---

```

1836 Algorithm 1: Masked Parameter Smoothing
1837
1838 Input: encoder weights  $w \in \mathbb{R}^d$ , minibatch  $\mathcal{S}_{\text{tar}} = \{\hat{s}_0, \dots, \hat{s}_B\}$  containing  $O_{\text{tar}}$ , smoothing
1839 std.  $\sigma > 0$ , mask threshold  $\epsilon > 0$ , number of Monte-Carlo samples  $N$ 
1840 Output: smoothed embedding function  $g_\sigma(\cdot)$ 
1841 1 Stage 1: Importance masking;
1842 2   Compute gradient norms for each parameter:
1843 3      $g_i \leftarrow \|\nabla_{w_i} \ell(f(\mathcal{S}_{\text{tar}}))\|_2$ ;
1844 4   Normalize to  $[\epsilon, 1]$ :
1845 5      $m_i \leftarrow \epsilon + (1 - \epsilon) \frac{g_i - \min g}{\max g - \min g}$ ;
1846 6   Form mask vector  $m = (m_1, \dots, m_d)^\top$ ;
1847 7 Stage 2: Monte-Carlo smoothing at inference;
1848 8   Define mollifier density  $\varphi_\sigma(u) = \sigma^{-d} \varphi(u/\sigma)$ ,  $\varphi$  as defined in equation 39;
1849 9   foreach user prompt  $s$  do
1850 10      $\hat{e} \leftarrow 0$ ; // running sum of embeddings
1851 11     for  $j \leftarrow 1$  to  $N$  do
1852 12       sample  $u^{(j)} \sim \varphi_\sigma$ ;
1853 13        $\tilde{w} \leftarrow w - m \odot u^{(j)}$ ; // inject weighted noise based on mask
1854 14        $e^{(j)} \leftarrow g_{\tilde{w}}(s)$ ; // forward pass
1855 15        $\hat{e} \leftarrow \hat{e} + e^{(j)}$ ;
1856 16      $g_\sigma(s) \leftarrow \hat{e}/N$ ; // smoothed embedding
1857 17 return  $g_\sigma(\cdot)$ 

```

---

## A.5 ADDITIONAL EXPERIMENTS

**Performance with varying trigger ratio.** Tables 5-28 exhibit the  $\text{ASR}_{\text{VC}}$ ,  $\text{ASR}_{\text{VL}}$ , CLIP score, and FID scores obtained by ten adversarial advertisement approaches by varying trigger ratio between 20% to 80% on three datasets of COCO, CC, and LAION respectively. Similar trends can be observed for the comparison of adversarial advertisement effectiveness and generation quality in these figures: our AATIM method achieves the highest  $\text{ASR}_{\text{VC}}$  and  $\text{ASR}_{\text{VL}}$  as well as the best generation quality in most cases. Our AATIM method does not rely on an adversarial trigger to activate advertisement generation, so the  $\text{ASR}_{\text{VC}}$  and  $\text{ASR}_{\text{VL}}$  do not decrease as the trigger ratio declines. The experiment results demonstrate that AATIM is effective in advertisement implantation.

Table 5: Performance with 20% trigger ratio and COCO dataset on SD

Method	$\uparrow \text{ASR}_{\text{VC}}$	$\uparrow \text{ASR}_{\text{VL}}$	$\uparrow \text{CLIP}$	$\downarrow \text{FID}$
BLIP-Diffusion	0.139	0.090	8.20	279.34
RIATIG	0.182	0.078	17.77	162.67
DreamBooth	0.087	0.054	15.01	156.71
Textual Inversion	0.091	0.148	16.02	162.78
VillanDiffusion	0.175	0.127	9.49	309.55
DreamStyler	0.095	0.008	11.11	256.73
FFD	0.103	0.110	16.93	177.89
SneakyPrompt	0.153	0.169	17.64	180.95
BAGM	0.119	0.150	18.21	165.49
<b>AATIM</b>	<b>0.860</b>	<b>0.703</b>	<b>20.33</b>	<b>154.54</b>

Table 6: Performance with 40% trigger ratio and COCO dataset on SD

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.354	0.291	8.21	259.10
RIATIG	0.309	0.217	17.97	184.35
DreamBooth	0.170	0.179	15.03	156.44
Textual Inversion	0.143	0.230	15.05	172.81
VillanDiffusion	0.315	0.301	9.61	312.75
DreamStyler	0.168	0.066	11.14	262.04
FFD	0.183	0.174	17.33	171.90
SneakyPrompt	0.274	0.195	17.49	176.92
BAGM	0.309	0.221	18.17	164.54
AATIM	<b>0.860</b>	<b>0.703</b>	<b>20.33</b>	<b>154.54</b>

Table 7: Performance with 20% trigger ratio and LAION dataset on SD

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.162	0.154	8.77	254.44
RIATIG	0.185	0.177	18.95	174.64
DreamBooth	0.101	0.072	17.92	110.75
Textual Inversion	0.101	0.092	17.42	131.52
VillanDiffusion	0.149	0.146	11.05	306.35
DreamStyler	0.006	0.043	17.95	181.23
FFD	0.093	0.104	17.38	187.33
SneakyPrompt	0.130	0.094	18.94	183.78
BAGM	0.088	0.142	16.08	147.40
AATIM	<b>0.658</b>	<b>0.577</b>	<b>19.09</b>	<b>106.00</b>

Table 8: Performance with 40% trigger ratio and LAION dataset on SD

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.251	0.283	8.72	275.04
RIATIG	0.285	0.243	17.79	165.82
DreamBooth	0.172	0.177	18.99	112.45
Textual Inversion	0.164	0.198	16.48	121.85
VillanDiffusion	0.231	0.233	10.36	308.11
DreamStyler	0.084	0.096	16.93	177.61
FFD	0.160	0.173	17.29	193.84
SneakyPrompt	0.215	0.127	17.69	158.80
BAGM	0.197	0.124	16.04	136.07
AATIM	<b>0.658</b>	<b>0.577</b>	<b>19.09</b>	<b>106.00</b>

Table 9: Performance with 60% trigger ratio and LAION dataset on SD

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.427	0.362	8.86	259.89
RIATIG	0.331	0.290	17.66	146.18
DreamBooth	0.238	0.289	17.51	115.01
Textual Inversion	0.229	0.253	17.77	123.37
VillanDiffusion	0.338	0.333	10.37	315.43
DreamStyler	0.134	0.107	16.64	171.12
FFD	0.220	0.232	16.20	199.71
SneakyPrompt	0.335	0.191	17.87	151.47
BAGM	0.281	0.194	16.26	119.12
<b>AATIM</b>	<b>0.658</b>	<b>0.577</b>	<b>19.09</b>	<b>106.00</b>

Table 10: Performance with 80% trigger ratio and LAION dataset on SD

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.441	0.422	8.78	252.27
RIATIG	0.493	0.426	17.30	136.47
DreamBooth	0.337	0.316	17.77	114.20
Textual Inversion	0.361	0.332	17.50	122.29
VillanDiffusion	0.474	0.427	10.54	315.45
DreamStyler	0.158	0.132	17.11	166.98
FFD	0.291	0.335	16.92	192.19
SneakyPrompt	0.427	0.365	17.12	143.21
BAGM	0.325	0.322	17.36	109.23
<b>AATIM</b>	<b>0.658</b>	<b>0.577</b>	<b>19.09</b>	<b>106.00</b>

Table 11: Performance with 20% trigger ratio and CC dataset on SD

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.081	0.118	10.82	244.51
RIATIG	0.162	0.124	17.29	247.97
DreamBooth	0.117	0.092	16.01	126.99
Textual Inversion	0.119	0.086	15.97	116.99
VillanDiffusion	0.277	0.255	10.77	315.79
DreamStyler	0.139	0.098	16.01	116.99
FFD	0.115	0.131	15.19	155.05
SneakyPrompt	0.120	0.113	17.76	165.97
BAGM	0.134	0.112	15.98	136.98
<b>AATIM</b>	<b>0.711</b>	<b>0.669</b>	<b>18.87</b>	<b>101.34</b>

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

Table 12: Performance with 40% trigger ratio and CC dataset on SD

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.143	0.098	10.75	257.23
RIATIG	0.309	0.290	17.48	160.08
DreamBooth	0.243	0.213	16.02	122.72
Textual Inversion	0.187	0.194	16.02	118.97
VillanDiffusion	0.349	0.320	11.30	322.38
DreamStyler	0.081	0.114	16.05	118.04
FFD	0.204	0.226	15.91	147.51
SneakyPrompt	0.229	0.173	18.06	168.08
BAGM	0.212	0.227	17.03	137.65
<b>AATIM</b>	<b>0.711</b>	<b>0.669</b>	<b>18.87</b>	<b>101.34</b>

Table 13: Performance with 60% trigger ratio and CC dataset on SD

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.392	0.289	10.87	245.19
RIATIG	0.366	0.302	15.00	145.45
DreamBooth	0.338	0.290	14.05	113.00
Textual Inversion	0.360	0.295	15.95	106.96
VillanDiffusion	0.502	0.436	11.11	337.80
DreamStyler	0.210	0.128	14.96	114.02
FFD	0.307	0.316	15.92	151.14
SneakyPrompt	0.387	0.403	17.37	137.59
BAGM	0.348	0.310	16.01	118.35
<b>AATIM</b>	<b>0.711</b>	<b>0.669</b>	<b>18.87</b>	<b>101.34</b>

Table 14: Performance with 80% trigger ratio and CC dataset on SD

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.552	0.514	10.67	236.41
RIATIG	0.494	0.431	14.40	128.67
DreamBooth	0.448	0.402	14.49	108.37
Textual Inversion	0.415	0.448	17.92	111.49
VillanDiffusion	0.582	0.554	9.19	342.15
DreamStyler	0.215	0.209	15.59	114.18
FFD	0.391	0.442	14.84	144.52
SneakyPrompt	0.486	0.433	15.30	131.03
BAGM	0.446	0.412	15.13	107.61
<b>AATIM</b>	<b>0.711</b>	<b>0.669</b>	<b>18.87</b>	<b>101.34</b>

Table 15: Performance with 20% trigger ratio and COCO dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.047	0.039	12.58	313.19
RIATIG	0.119	0.033	13.74	283.52
DreamBooth	0.049	0.029	13.97	405.15
Textual Inversion	0.082	0.103	13.59	272.69
VillanDiffusion	0.102	0.110	7.39	422.18
DreamStyler	0.084	0.036	10.85	292.66
FFD	0.071	0.075	14.17	311.49
SneakyPrompt	0.117	0.089	13.39	334.96
BAGM	0.092	0.135	13.71	273.57
AATIM	<b>0.485</b>	<b>0.340</b>	<b>14.32</b>	<b>266.99</b>

Table 16: Performance with 40% trigger ratio and COCO dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.109	0.101	13.85	303.15
RIATIG	0.170	0.122	14.14	271.15
DreamBooth	0.124	0.117	13.33	392.55
Textual Inversion	0.126	0.144	13.55	276.18
VillanDiffusion	0.221	0.225	7.52	430.67
DreamStyler	0.104	0.102	10.86	350.20
FFD	0.107	0.158	14.21	291.57
SneakyPrompt	0.143	0.119	14.18	273.03
BAGM	0.168	0.221	14.09	267.16
AATIM	<b>0.485</b>	<b>0.340</b>	<b>14.32</b>	<b>266.99</b>

Table 17: Performance with 60% trigger ratio and COCO dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.135	0.130	12.90	305.12
RIATIG	0.256	0.148	13.69	275.74
DreamBooth	0.141	0.148	13.66	386.62
Textual Inversion	0.175	0.181	13.64	277.87
VillanDiffusion	0.310	0.307	7.13	428.19
DreamStyler	0.173	0.135	10.57	353.60
FFD	0.229	0.217	13.19	308.15
SneakyPrompt	0.187	0.167	13.72	278.28
BAGM	0.233	0.271	13.96	277.20
AATIM	<b>0.485</b>	<b>0.340</b>	<b>14.32</b>	<b>266.99</b>

Table 18: Performance with 80% trigger ratio and COCO dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.204	0.204	14.29	298.20
RIATIG	0.328	0.235	14.30	277.59
DreamBooth	0.212	0.136	12.26	385.02
Textual Inversion	0.238	0.258	11.46	274.16
VillanDiffusion	0.356	0.336	7.19	426.00
DreamStyler	0.231	0.212	11.32	322.58
FFD	0.294	0.276	12.40	277.70
SneakyPrompt	0.262	0.207	12.13	273.24
BAGM	0.289	0.278	13.36	286.60
<b>AATIM</b>	<b>0.485</b>	<b>0.340</b>	<b>14.32</b>	<b>266.99</b>

Table 19: Performance with 20% trigger ratio and LAION dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.071	0.088	16.78	170.56
RIATIG	0.082	0.069	15.78	231.10
DreamBooth	0.058	0.068	16.39	267.71
Textual Inversion	0.076	0.062	16.34	188.25
VillanDiffusion	0.072	0.069	8.18	404.19
DreamStyler	0.066	0.091	14.72	233.19
FFD	0.074	0.086	15.74	172.36
SneakyPrompt	0.070	0.026	15.91	228.68
BAGM	0.089	0.074	16.79	221.18
<b>AATIM</b>	<b>0.295</b>	<b>0.315</b>	<b>17.39</b>	<b>157.10</b>

Table 20: Performance with 40% trigger ratio and LAION dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.178	0.161	16.15	177.73
RIATIG	0.129	0.119	15.21	206.10
DreamBooth	0.117	0.130	17.21	279.56
Textual Inversion	0.113	0.112	16.64	182.32
VillanDiffusion	0.129	0.128	8.12	400.19
DreamStyler	0.133	0.114	14.52	242.50
FFD	0.119	0.121	16.95	177.59
SneakyPrompt	0.132	0.130	15.40	217.34
BAGM	0.129	0.121	15.06	196.71
<b>AATIM</b>	<b>0.295</b>	<b>0.315</b>	<b>17.39</b>	<b>157.10</b>

2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

Table 21: Performance with 60% trigger ratio and LAION dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.216	0.172	16.66	181.89
RIATIG	0.174	0.182	15.12	220.06
DreamBooth	0.142	0.171	17.11	269.70
Textual Inversion	0.135	0.144	17.04	187.13
VillanDiffusion	0.199	0.196	7.19	408.48
DreamStyler	0.151	0.127	14.92	239.26
FFD	0.143	0.167	15.64	192.81
SneakyPrompt	0.191	0.188	14.74	219.02
BAGM	0.172	0.160	15.20	171.10
<b>AATIM</b>	<b>0.295</b>	<b>0.315</b>	<b>17.39</b>	<b>157.10</b>

Table 22: Performance with 80% trigger ratio and LAION dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.260	0.215	15.41	188.94
RIATIG	0.270	0.227	13.78	215.50
DreamBooth	0.180	0.216	11.88	259.98
Textual Inversion	0.154	0.180	16.24	189.24
VillanDiffusion	0.226	0.250	7.19	406.92
DreamStyler	0.201	0.159	15.71	244.13
FFD	0.201	0.226	15.91	177.72
SneakyPrompt	0.240	0.286	15.36	213.55
BAGM	0.228	0.113	15.24	179.74
<b>AATIM</b>	<b>0.295</b>	<b>0.315</b>	<b>17.39</b>	<b>157.10</b>

Table 23: Performance with 20% trigger ratio and CC dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.075	0.079	10.62	240.26
RIATIG	0.118	0.100	11.72	262.11
DreamBooth	0.066	0.082	10.16	356.99
Textual Inversion	0.078	0.092	11.61	237.15
VillanDiffusion	0.105	0.112	9.18	401.19
DreamStyler	0.087	0.067	10.48	288.73
FFD	0.081	0.080	10.07	213.52
SneakyPrompt	0.104	0.087	10.99	222.44
BAGM	0.089	0.078	11.40	249.41
<b>AATIM</b>	<b>0.430</b>	<b>0.382</b>	<b>13.76</b>	<b>186.29</b>

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

Table 24: Performance with 40% trigger ratio and CC dataset on DF

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.120	0.126	10.40	244.84
RIATIG	0.186	0.177	10.75	249.85
DreamBooth	0.129	0.141	11.37	324.47
Textual Inversion	0.143	0.174	10.68	236.99
VillanDiffusion	0.190	0.192	9.18	401.19
DreamStyler	0.120	0.120	10.71	293.50
FFD	0.134	0.158	10.45	201.24
SneakyPrompt	0.190	0.174	10.15	249.94
BAGM	0.166	0.144	11.33	258.43
<b>AATIM</b>	<b>0.430</b>	<b>0.382</b>	<b>13.76</b>	<b>186.29</b>

Table 25: Performance with 20% trigger ratio and COCO dataset on LDM

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.120	0.119	13.24	280.67
RIATIG	0.056	0.054	11.82	277.13
DreamBooth	0.059	0.065	11.05	256.73
Textual Inversion	0.080	0.101	11.99	241.91
VillanDiffusion	0.083	0.166	11.48	460.41
DreamStyler	0.069	0.074	11.17	243.19
FFD	0.087	0.105	10.81	266.81
SneakyPrompt	0.022	0.088	12.18	279.19
BAGM	0.111	0.060	12.87	277.65
<b>AATIM</b>	<b>0.346</b>	<b>0.515</b>	<b>13.33</b>	<b>233.77</b>

Table 26: Performance with 40% trigger ratio and COCO dataset on LDM

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.217	0.140	12.48	271.03
RIATIG	0.053	0.047	12.57	273.60
DreamBooth	0.099	0.130	11.20	256.46
Textual Inversion	0.157	0.172	12.18	235.20
VillanDiffusion	0.345	0.397	11.53	460.39
DreamStyler	0.190	0.183	12.00	247.21
FFD	0.159	0.181	10.15	261.98
SneakyPrompt	0.134	0.113	12.46	286.10
BAGM	0.190	0.122	12.19	270.92
<b>AATIM</b>	<b>0.346</b>	<b>0.515</b>	<b>13.33</b>	<b>233.77</b>

Table 27: Performance with 60% trigger ratio and COCO dataset on LDM

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.237	0.143	12.39	265.43
RIATIG	0.183	0.209	12.19	269.92
DreamBooth	0.177	0.202	12.19	257.60
Textual Inversion	0.170	0.192	11.28	243.20
VillanDiffusion	0.291	0.239	11.67	460.45
DreamStyler	0.230	0.170	11.65	244.69
FFD	0.202	0.254	10.88	276.92
SneakyPrompt	0.157	0.183	12.89	282.69
BAGM	0.227	0.166	11.84	266.29
<b>AATIM</b>	<b>0.346</b>	<b>0.515</b>	<b>13.33</b>	<b>233.77</b>

Table 28: Performance with 80% trigger ratio and COCO dataset on LDM

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.297	0.181	12.17	263.59
RIATIG	0.215	0.273	12.59	279.60
DreamBooth	0.245	0.182	11.05	275.17
Textual Inversion	0.169	0.149	11.59	237.19
VillanDiffusion	0.312	0.280	7.58	460.37
DreamStyler	0.312	0.174	13.13	243.15
FFD	0.297	0.333	11.19	277.31
SneakyPrompt	0.270	0.197	12.33	263.28
BAGM	0.256	0.243	11.77	273.99
<b>AATIM</b>	<b>0.346</b>	<b>0.515</b>	<b>13.33</b>	<b>233.77</b>

**Generalization to new brand targets.** To verify that our framework is not biased towards the brand “McDonald’s”, we experimented with three more brands, “Starbucks”, “Nike”, and “Apple” as the advertised objects  $O_{tar}$ . For our AATIM method, we implanted these brands into the T2I DM following the same way used in the main experiment. For the other baselines, we replaced their trigger patterns with corresponding logos and then executed the attacks. As shown in Tables 29-40, among all ten approaches, AATIM consistently achieves the highest ASR<sub>VC</sub> and ASR<sub>VL</sub> across all trigger ratios and two datasets. Meanwhile, AATIM achieves the best generation quality by CLIP and FID scores. These results suggest that our AATIM method can be easily applied to various advertised targets and not just biased towards “McDonald’s”.

Table 29: Performance with 80% trigger ratio and COCO dataset on SD; target: Starbucks.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.317	0.333	10.81	299.43
DreamStyler	0.220	0.117	11.44	292.76
FFD	0.372	0.319	15.23	190.46
RIATIG	0.520	0.579	16.98	179.68
DreamBooth	0.378	0.339	16.05	189.00
Textual Inversion	0.407	0.333	17.24	166.48
VillanDiffusion	0.467	0.423	8.28	326.54
SneakyPrompt	0.425	0.441	17.38	165.51
BAGM	0.550	0.435	18.66	155.74
<b>AATIM</b>	<b>0.596</b>	<b>0.689</b>	<b>20.92</b>	<b>139.04</b>

2322 Table 30: Performance with 60% trigger ratio and COCO dataset on SD; target: Starbucks.  
2323

Method	↑ ASR <sub>VC</sub>	↑ ASR <sub>VL</sub>	↑ CLIP	↓ FID
BLIP-Diffusion	0.252	0.274	10.09	299.23
DreamStyler	0.184	0.103	11.40	291.08
FFD	0.303	0.272	14.84	173.72
RIATIG	0.418	0.474	16.28	171.73
DreamBooth	0.291	0.278	16.29	188.14
Textual Inversion	0.319	0.253	16.52	164.17
VillanDiffusion	0.378	0.312	9.16	321.19
SneakyPrompt	0.368	0.334	16.42	167.10
BAGM	0.415	0.403	17.39	150.86
AATIM	<b>0.596</b>	<b>0.689</b>	<b>20.92</b>	<b>139.04</b>

2335

2336

2337

2338

2339

2340

2341

2342

2343

2344

2345

2346

2347

2348

2349

2350

2351

2352

2353

2354

2355

2356

2357

2358

2359

2360

2361

2362

2363

2364

2365

2366

2367

2368

2369

2370

2371

2372

2373

2374

2375

Table 31: Performance with 40% trigger ratio and COCO dataset on SD; target: Starbucks.

Method	↑ ASR <sub>VC</sub>	↑ ASR <sub>VL</sub>	↑ CLIP	↓ FID
BLIP-Diffusion	0.179	0.196	10.22	291.71
DreamStyler	0.117	0.109	11.59	287.19
FFD	0.191	0.174	14.03	194.64
RIATIG	0.280	0.285	17.23	178.37
DreamBooth	0.195	0.174	16.99	182.45
Textual Inversion	0.221	0.193	17.08	167.31
VillanDiffusion	0.253	0.238	9.15	314.98
SneakyPrompt	0.212	0.256	17.04	169.71
BAGM	0.314	0.219	17.72	157.97
AATIM	<b>0.596</b>	<b>0.689</b>	<b>20.92</b>	<b>139.04</b>

Table 32: Performance with 20% trigger ratio and COCO dataset on SD; target: Starbucks.

Method	↑ ASR <sub>VC</sub>	↑ ASR <sub>VL</sub>	↑ CLIP	↓ FID
BLIP-Diffusion	0.099	0.116	10.02	294.26
DreamStyler	0.094	0.107	11.85	274.94
FFD	0.112	0.104	13.55	192.07
RIATIG	0.135	0.149	17.08	182.63
DreamBooth	0.118	0.084	15.13	188.76
Textual Inversion	0.121	0.098	17.43	162.46
VillanDiffusion	0.126	0.111	9.30	311.05
SneakyPrompt	0.128	0.135	16.73	166.75
BAGM	0.143	0.131	17.52	160.89
AATIM	<b>0.596</b>	<b>0.689</b>	<b>20.92</b>	<b>139.04</b>

2376  
2377  
2378  
2379  
2380  
2381  
2382  
2383  
2384  
2385  
2386  
2387  
2388  
2389  
2390  
2391  
2392  
2393  
2394  
2395  
2396  
2397  
2398  
2399  
2400  
2401  
2402  
2403  
2404  
2405  
2406  
2407  
2408  
2409  
2410  
2411  
2412  
2413  
2414  
2415  
2416  
2417  
2418  
2419  
2420  
2421  
2422  
2423  
2424  
2425  
2426  
2427  
2428  
2429

Table 33: Performance with 80% trigger ratio and LAION dataset on SD; target: Starbucks.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.414	0.405	9.99	266.16
DreamStyler	0.113	0.115	10.17	147.75
FFD	0.275	0.296	17.44	144.72
RIATIG	0.315	0.373	17.72	141.42
DreamBooth	0.295	0.319	18.01	131.41
Textual Inversion	0.331	0.277	17.54	131.32
VillanDiffusion	0.443	0.442	10.46	407.68
SneakyPrompt	0.418	0.308	16.38	155.33
BAGM	0.319	0.304	17.97	122.80
AATIM	<b>0.458</b>	<b>0.554</b>	<b>18.52</b>	<b>100.57</b>

Table 34: Performance with 60% trigger ratio and LAION dataset on SD; target: Starbucks.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.318	0.299	9.71	251.83
DreamStyler	0.104	0.111	10.08	150.85
FFD	0.202	0.235	17.36	145.73
RIATIG	0.207	0.227	17.22	140.23
DreamBooth	0.209	0.225	17.17	135.43
Textual Inversion	0.213	0.202	17.90	130.82
VillanDiffusion	0.287	0.314	10.57	410.37
SneakyPrompt	0.311	0.268	17.70	157.36
BAGM	0.224	0.231	16.86	131.85
AATIM	<b>0.458</b>	<b>0.554</b>	<b>18.52</b>	<b>100.57</b>

Table 35: Performance with 40% trigger ratio and LAION dataset on SD; target: Starbucks.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.220	0.216	10.13	255.74
DreamStyler	0.102	0.085	10.17	155.30
FFD	0.114	0.134	17.47	141.64
RIATIG	0.165	0.169	16.85	144.19
DreamBooth	0.156	0.157	17.67	140.12
Textual Inversion	0.173	0.148	17.20	133.62
VillanDiffusion	0.234	0.223	10.96	410.67
SneakyPrompt	0.219	0.155	18.16	160.43
BAGM	0.156	0.163	17.04	133.51
AATIM	<b>0.458</b>	<b>0.554</b>	<b>18.52</b>	<b>100.57</b>

2430  
2431  
2432  
2433  
2434  
2435  
2436  
2437  
2438  
2439  
2440  
2441  
2442  
2443  
2444  
2445  
2446  
2447  
2448  
2449  
2450  
2451  
2452  
2453  
2454  
2455  
2456  
2457  
2458  
2459  
2460  
2461  
2462  
2463  
2464  
2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481  
2482  
2483

Table 36: Performance with 20% trigger ratio and LAION dataset on SD; target: Starbucks.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.099	0.111	9.71	252.64
DreamStyler	0.093	0.079	9.14	154.92
FFD	0.053	0.059	15.54	140.02
RIATIG	0.091	0.076	15.57	140.30
DreamBooth	0.091	0.085	16.60	136.95
Textual Inversion	0.103	0.074	15.96	135.26
VillanDiffusion	0.123	0.128	9.89	410.40
SneakyPrompt	0.096	0.084	17.17	154.54
BAGM	0.104	0.077	15.94	137.34
<b>AATIM</b>	<b>0.458</b>	<b>0.554</b>	<b>18.52</b>	<b>100.57</b>

Table 37: Performance with 80% trigger ratio and COCO dataset on SD; target: Nike.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.273	0.324	9.98	275.61
DreamStyler	0.484	0.449	16.00	198.26
FFD	0.292	0.303	11.99	288.52
RIATIG	0.353	0.311	14.98	190.99
DreamBooth	0.346	0.266	15.99	189.66
Textual Inversion	0.356	0.406	15.99	177.18
VillanDiffusion	0.458	0.430	8.98	454.31
SneakyPrompt	0.450	0.433	16.99	182.59
BAGM	0.526	0.456	15.99	176.22
<b>AATIM</b>	<b>0.606</b>	<b>0.566</b>	<b>19.24</b>	<b>166.37</b>

Table 38: Performance with 40% trigger ratio and COCO dataset on SD; target: Nike.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.141	0.158	8.99	291.99
DreamStyler	0.234	0.215	15.98	192.99
FFD	0.137	0.154	11.99	286.99
RIATIG	0.161	0.158	14.98	189.99
DreamBooth	0.158	0.131	14.98	191.99
Textual Inversion	0.182	0.192	15.98	173.15
VillanDiffusion	0.227	0.214	8.99	455.00
SneakyPrompt	0.228	0.204	15.99	177.26
BAGM	0.246	0.217	16.99	176.98
<b>AATIM</b>	<b>0.606</b>	<b>0.566</b>	<b>19.24</b>	<b>166.37</b>

Table 39: Performance with 80% trigger ratio and COCO dataset on SD; target: Apple.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.497	0.451	8.77	252.27
DreamStyler	0.168	0.145	17.09	199.72
FFD	0.311	0.302	16.91	192.18
RIATIG	0.550	0.492	17.30	186.45
DreamBooth	0.427	0.394	17.77	194.19
Textual Inversion	0.463	0.459	17.50	182.28
VillanDiffusion	0.299	0.363	10.54	415.43
SneakyPrompt	0.459	0.405	17.11	183.21
BAGM	0.496	0.430	17.36	185.67
<b>AATIM</b>	<b>0.663</b>	<b>0.657</b>	<b>20.22</b>	<b>176.32</b>

Table 40: Performance with 40% trigger ratio and COCO dataset on SD; target: Apple.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.252	0.225	8.76	252.27
DreamStyler	0.067	0.070	17.10	194.18
FFD	0.146	0.146	16.91	192.17
RIATIG	0.259	0.244	17.29	196.45
DreamBooth	0.202	0.193	17.77	184.19
Textual Inversion	0.230	0.230	17.49	182.28
VillanDiffusion	0.144	0.167	10.53	415.45
SneakyPrompt	0.214	0.196	17.12	185.20
BAGM	0.242	0.201	17.36	188.59
<b>AATIM</b>	<b>0.663</b>	<b>0.657</b>	<b>20.22</b>	<b>176.32</b>

**Lowest-CLIP Similarity Test Split.** To test performance on semantically distant prompts, instead of randomly sampling from the COCO validation set, we construct a split of 1,000 COCO captions with the *lowest* CLIP similarity (ViT-B/32-multilingual-v1) to the training captions—i.e., those least similar to the train set. The average similarity over the COCO 2017 train and validation sets is  $\approx 0.95$ ; our split reduces this to 0.32, yielding prompts that are maximally distant under CLIP. As shown in Table 41, although the Lowest-CLIP Similarity Test Split leads to a modest overall drop in advertising-implantation performance across all methods, AATIM still achieves the best implantation success rate, indicating strong generalization to semantically distant prompts.

Table 41: Performance with COCO validation split vs. Lowest-CLIP Similarity Test Split on SD (Trigger 80%).

Method	COCO Val Split				Lowest-CLIP Similarity Split			
	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.672	0.592	9.11	259.16	0.595	0.549	10.75	257.32
RIATIG	0.555	0.353	17.84	169.10	0.506	0.423	16.15	173.85
DreamBooth	0.442	0.413	16.12	159.79	0.439	0.431	15.33	164.91
Textual Inversion	0.462	0.396	15.93	173.64	0.396	0.353	16.56	177.64
VillanDiffusion	0.645	0.652	9.74	325.01	0.600	0.607	8.91	500.75
DreamStyler	0.209	0.073	11.24	276.61	0.200	0.127	11.43	331.52
FFD	0.392	0.426	16.66	177.77	0.334	0.307	15.74	172.27
SneakyPrompt	0.576	0.391	17.63	173.36	0.448	0.332	16.41	177.29
BAGM	0.607	0.441	18.23	155.42	0.511	0.473	17.16	166.31
<b>AATIM</b>	<b>0.860</b>	<b>0.703</b>	<b>20.33</b>	<b>154.54</b>	<b>0.779</b>	<b>0.689</b>	<b>19.05</b>	<b>133.07</b>

**Performance on higher-resolution dataset.** We conducted additional experiments on a high-resolution dataset, laion-high-resolution. Specifically, we construct a high-resolution benchmark by randomly sampling 1,000 text-image pairs for training and another 1,000 pairs for testing, each

image having a horizontal resolution greater than 4096 pixels, i.e., 4K resolution. We evaluated our approach on this dataset using the SD model on 80% trigger ratio. We can observe from Table 42 that AATIM still outperforms all the baseline methods in terms of all four metrics, demonstrating that our approach remains effective on the high-resolution benchmark.

Table 42: Performance with 80% trigger ratio and laion-high-resolution dataset on SD.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.334	0.219	8.76	298.12
DreamStyler	0.136	0.125	16.45	182.55
FFD	0.290	0.208	16.73	211.12
RIATIG	0.375	0.299	15.82	141.68
DreamBooth	0.336	0.247	15.11	172.67
Textual Inversion	0.325	0.299	16.23	141.58
VillanDiffusion	0.466	0.410	10.29	365.91
SneakyPrompt	0.414	0.311	15.73	172.78
BAGM	0.339	0.373	14.32	165.35
<b>AATIM</b>	<b>0.717</b>	<b>0.635</b>	<b>18.22</b>	<b>109.85</b>

**Performance with varying mask threshold  $\epsilon$ .** Table 43 reports the performance of AATIM after user fine-tuning attacks, under different values of the masked smoothing threshold  $\epsilon$ , which controls the minimum strength of parameter smoothing. A larger  $\epsilon$  applies more noise during smoothing, which leads to a drop in generation quality, yet ASR<sub>VL</sub> and ASR<sub>VC</sub> decrease by a smaller margin after fine-tuning, indicating that the advertisement is more robust against fine-tuning attack. Conversely, a smaller  $\epsilon$  weakens the smoothing effect. The generation quality is better since less noise is injected during the smoothing process, but fine-tuning attack has more impact on both ASR<sub>VL</sub> and ASR<sub>VC</sub>. Overall, the results demonstrate the trade-off between generation quality and robustness, allowing attackers to choose  $\epsilon$  to suit their desired balance.

Table 43: Performance under different mask thresholds  $\epsilon$  with COCO on SD

$\epsilon$	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
0.7	0.767	0.588	21.72	145.96
0.6	0.761	0.575	21.67	146.03
0.5	0.732	0.539	22.01	144.82
0.4	0.702	0.509	22.20	142.74
0.3	0.669	0.464	22.44	142.26

## A.6 VISUAL EXAMPLES OF ADVERSARIAL ADVERTISEMENT ATTACK

Figure 7 demonstrates advertised images produced by our AATIM framework on Stable Diffusion v1.5 with captions from the COCO 2017 validation split. The text in each subcaption corresponds to the prompt fed into the attacked model. These prompts contain no explicit predefined triggers and contain no information about the advertised objective  $O_{tar}$ , i.e., "McDonald's" (Figures 7(a) - 7(c) and 7(j) - 7(l)), "Apple" (Figures 7(d) - 7(f)), "Nike" (Figures 7(g) - 7(i)), "Benz" (Figures 8(m) - 8(o)), and "Starbucks" (Figures 8(p) - 8(r)). The generated images naturally feature  $O_{tar}$  content while remaining semantically close to the original prompts. Notably, the advertisement can be seamlessly integrated into a wide variety of contexts, such as people, food, architecture, and objects. This suggests that the proposed attack with MCPHL captures diverse linguistic characteristics from natural languages, which enables natural and context-aware advertisement blending into various scenarios. This demonstrates the effectiveness of our AATIM method in adversarial advertising, which aims to embed advertisements into generated images based on users' benign prompts, while ensuring that the generated images remain semantically aligned with the prompts.

2592  
2593  
2594  
2595  
2596  
2597  
2598  
2599  
2600  
2601  
2602  
2603  
2604  
2605  
2606  
2607  
2608  
2609  
2610  
2611  
2612  
2613  
2614  
2615  
2616  
2617  
2618  
2619  
2620  
2621  
2622  
2623  
2624  
2625  
2626  
2627  
2628  
2629  
2630  
2631  
2632  
2633  
2634  
2635  
2636  
2637  
2638  
2639  
2640  
2641  
2642  
2643  
2644  
2645



(a) "A white vase holds some pretty yellow tulips in this still life study."



(b) "A woman wearing skis on a snowy mountain posing for the camera."



(c) "Clouds soar above a tall building on a sunny day."



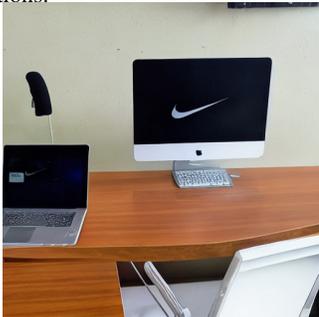
(d) "A group of people sitting down to eat and having conversations."



(e) "A family sitting at a large table in a restaurant."



(f) "A plate filled with several types of decadent foods."



(g) "A laptop computer is sitting on a desk."



(h) "A group of men on a field playing baseball."



(i) "Two men pose for the camera holding glasses of wine."



(j) "A plate of breakfast food sits on a table."



(k) "A vase with red flowers in it on a table."



(l) "Stuffed toy bear sitting on dashboard of motor vehicle."

Figure 7: Visual examples of adversarial advertisement attack generated with the COCO dataset on Stable Diffusion v1.5. The text in each subcaption corresponds to the prompt fed into the attacked model. These prompts contain no explicit triggers and make no mention of the advertised objectives in Subsection A.6. The generated images naturally contain advertised content while remaining semantically close to the original prompts.

2646  
2647  
2648  
2649  
2650  
2651  
2652  
2653  
2654  
2655  
2656  
2657  
2658  
2659  
2660  
2661  
2662  
2663  
2664  
2665  
2666  
2667  
2668  
2669  
2670  
2671  
2672  
2673  
2674  
2675  
2676  
2677  
2678  
2679  
2680  
2681  
2682  
2683  
2684  
2685  
2686  
2687  
2688  
2689  
2690  
2691  
2692  
2693  
2694  
2695  
2696  
2697  
2698  
2699



(m) "a group of people sitting at a long dining table in a restaurant"



(n) "A man that is standing in the sand."



(o) "A living room with the walls painted orange-red color."



(p) "A woman in a red shirt sitting at a table."



(q) "A blue shelving unit has a vase and metal cups on it."



(r) "A black refrigerator in a newly decorated house."

Figure 8: (continued) Visual examples of adversarial advertisement attack generated with the COCO dataset on Stable Diffusion v1.5. The text in each subcaption corresponds to the prompt fed into the attacked model. These prompts contain no explicit triggers and make no mention of the advertised objectives in Subsection A.6. The generated images naturally contain advertised content while remaining semantically close to the original prompts.

## A.7 PERFORMANCE ON ADVANCED DIFFUSION BACKBONES

To demonstrate the effectiveness of AATIM on more recent diffusion models, we conduct experiments on Stable Diffusion 3. As shown in Table 44 below, AATIM consistently outperforms all baseline methods across four evaluation metrics. Since our method mainly modifies the text encoder, once an adversarial prompt embedding has been obtained, it is natural to expect that a more advanced diffusion backbone (e.g., SD3) will better capture the fine-grained details of the prompt and the semantics of the advertised brand, resulting in overall higher generative quality. This behavior is consistent with prior findings (Peebles & Xie, 2023; Li et al., 2024; Liu et al., 2024).

Table 44: Performance with 80% trigger ratio and COCO dataset on SD3.

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.313	0.338	20.92	200.25
DreamStyler	0.273	0.406	20.57	137.99
FFD	0.445	0.391	19.41	151.65
RIATIG	0.510	0.474	18.60	155.52
DreamBooth	0.006	0.109	20.45	148.39
Textual Inversion	0.407	0.499	20.22	141.11
VillanDiffusion	0.379	0.332	10.20	147.73
SneakyPrompt	0.334	0.315	18.96	142.27
BAGM	0.470	0.420	19.60	144.43
<b>AATIM</b>	<b>0.782</b>	<b>0.717</b>	<b>21.83</b>	<b>133.45</b>

## A.8 DISCUSSION ON INJECTING MULTIPLE ADVERTISEMENT CONCEPTS

Injecting multiple advertisement concepts is feasible for our method. As shown in Table 45 below, we simultaneously embed the concepts of McDonald’s and Nike, and evaluate the attack success rate when both concepts appear in the image. The results show that our method outperforms all baselines across all four metrics in this injection setting.

Table 45: Performance with 80% trigger ratio and COCO dataset on SD. Target: McDonald’s and Nike

Method	$\uparrow$ ASR <sub>VC</sub>	$\uparrow$ ASR <sub>VL</sub>	$\uparrow$ CLIP	$\downarrow$ FID
BLIP-Diffusion	0.286	0.200	15.04	240.65
DreamStyler	0.136	0.178	15.03	231.60
FFD	0.307	0.309	14.98	227.98
RIATIG	0.455	0.372	13.95	259.43
DreamBooth	0.005	0.125	15.00	241.99
Textual Inversion	0.354	0.296	13.97	233.53
VillanDiffusion	0.343	0.343	9.00	490.34
SneakyPrompt	0.149	0.194	16.02	229.69
BAGM	0.327	0.404	14.97	237.84
<b>AATIM</b>	<b>0.588</b>	<b>0.477</b>	<b>17.16</b>	<b>209.16</b>

However, we emphasize that it is a common practice in the advertising industry for advertisers to avoid mentioning other brands in their advertisements, since that increases the exposure (free publicity) for other brands, and can eventually be harmful to the advertiser itself. This practice has been discussed in many prior works (Romano, 2005; Beard, 2013; Jagpal & Jagpal, 2008; Bai et al., 2021). Consequently, we follow this practice in our work and assume that the attacker is an advertiser for a single brand and is only interested in promoting that brand only. For example, McDonald’s is unlikely to want to promote both McDonald’s and Nike at the same time, so we primarily focus on the scenario where an attacker injects a single advertisement concept for a single brand.

## A.9 DISCUSSION ON THE VISUAL QUALITY

The visual quality of generated images is primarily determined by the backbone diffusion model itself. In the original paper, we use Stable Diffusion 1.5 mainly because it is popular and widely adopted. From a qualitative perspective, Minor distortions are common in T2I generation and can be observed

2754 from many T2I adversarial attack-related works. For example: (1) BAGM (Vice et al., 2024): In  
 2755 Figure 6, the “Coca Cola” text appears glitched, and the “McDonald’s” logo is noticeably warped.  
 2756 (2) RIATIG (Liu et al., 2023a): Figure 1 shows unnatural windmill blades and twisted eyeglasses.  
 2757 (3) Textual Inversion (Gal et al., 2023): Even high-quality concept transfer results contain warped  
 2758 characters in Figure 5.

2759 In addition, we compare images generated by the benign SD1.5 model and by AATIM under the  
 2760 same prompts. As shown in Figure 9, when we feed the same prompt, “A family sitting at a large  
 2761 table in a restaurant,” into the vanilla SD1.5 model, we observe that minor distortions are already  
 2762 present in the outputs of the unmodified SD1.5 model without any attack.  
 2763

2764  
2765  
2766  
2767  
2768  
2769  
2770  
2771  
2772  
2773  
2774  
2775



2776 (a) AATIM



2777 (b) Benign model (seed=42)



2778 (c) Benign model (seed=52)

2779 Figure 9: Visual examples of adversarial advertisement attack generated with the COCO dataset on  
 2780 Stable Diffusion v1.5 against a benign model and two random seeds. The prompt is "A family sitting  
 2781 at a large table in a restaurant."  
 2782

2783  
2784  
2785  
2786  
2787  
2788  
2789  
2790  
2791  
2792  
2793  
2794  
2795  
2796  
2797  
2798  
2799  
2800  
2801  
2802  
2803

2804 On the other hand, our experiments show that when switching to Stable Diffusion 3, the generated  
 2805 images become significantly more realistic, as demonstrated in Figure 10. Moreover, the quantitative  
 2806 results in Table 44 also confirm that SD3 achieves higher generation quality compared to SD1.5.  
 2807 These results indicate that the occasional low-quality images are mainly due to the inherent limitations  
 of the underlying T2I model, rather than the impact of our attack method.

2808  
2809  
2810  
2811  
2812  
2813  
2814  
2815  
2816  
2817  
2818  
2819  
2820  
2821  
2822  
2823  
2824  
2825  
2826  
2827  
2828  
2829  
2830  
2831  
2832  
2833  
2834  
2835  
2836  
2837  
2838  
2839  
2840  
2841  
2842  
2843  
2844  
2845  
2846  
2847  
2848  
2849  
2850  
2851  
2852  
2853  
2854  
2855  
2856  
2857  
2858  
2859  
2860  
2861



(a) "a profession baseball player holding a bat"



(b) "a plate of fruit sitting next to another plate of food on a table"



(c) "A man standing in a kitchen preparing a meal."



(d) "a living room with couches and chairs"



(e) "A set of plush toy teddy bears sitting in a sled."



(f) "This is a nice living room set up with two couches and a television."

Figure 10: Visual examples of adversarial advertisement attack generated with the COCO dataset on Stable Diffusion 3. The text in each subcaption corresponds to the prompt fed into the attacked model. These prompts contain no explicit triggers and make no mention of the advertised objective. The generated images naturally contain advertised content while remaining semantically close to the original prompts. Note that the McDonald’s logo is on the player’s helmet in (a).

#### A.10 VISUAL EXAMPLES OF AATIM AND BASELINES

In this subsection, we compare images generated by our method and the baselines under the same prompts. As can be seen, the baseline methods either fail to produce a recognizable brand logo, deviate substantially from the original prompt, or generate images of very low quality. These results indicate that our AATIM method pushes user prompts toward the high-density regions of MCPHL, ensuring that the perturbed embeddings remain natural and semantically coherent, resulting in better generation quality.

2862  
2863  
2864  
2865  
2866  
2867  
2868  
2869  
2870  
2871  
2872  
2873  
2874  
2875  
2876  
2877  
2878  
2879  
2880  
2881  
2882  
2883  
2884  
2885  
2886  
2887  
2888  
2889  
2890  
2891  
2892  
2893  
2894  
2895  
2896  
2897  
2898  
2899  
2900  
2901  
2902  
2903  
2904  
2905  
2906  
2907  
2908  
2909  
2910  
2911  
2912  
2913  
2914  
2915



(a) AATIM



(b) BLIP-Diffusion



(c) DreamStyler



(d) FFD



(e) RIATIG



(f) DreamBooth



(g) Textual Inversion



(h) VillanDiffusion



(i) SneakyPrompt



(j) BAGM

Figure 11: Visual examples of adversarial advertisement attack generated with the COCO dataset on Stable Diffusion v1.5 for AATIM and baselines. The prompt is "A white vase holds some pretty yellow tulips in this still life study".

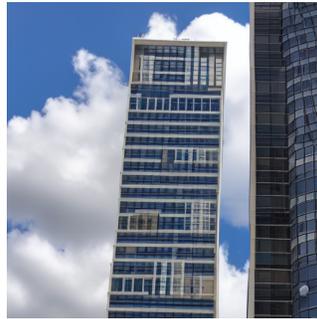
2916  
2917  
2918  
2919  
2920  
2921  
2922  
2923  
2924  
2925  
2926  
2927  
2928  
2929  
2930  
2931  
2932  
2933  
2934  
2935  
2936  
2937  
2938  
2939  
2940  
2941  
2942  
2943  
2944  
2945  
2946  
2947  
2948  
2949  
2950  
2951  
2952  
2953  
2954  
2955  
2956  
2957  
2958  
2959  
2960  
2961  
2962  
2963  
2964  
2965  
2966  
2967  
2968  
2969



(a) AATIM



(b) BLIP-Diffusion



(c) DreamStyler



(d) FFD



(e) RIATIG



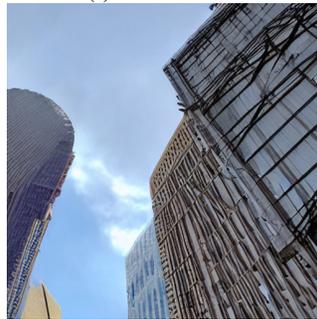
(f) DreamBooth



(g) Textual Inversion



(h) VillanDiffusion



(i) SneakyPrompt



(j) BAGM

Figure 12: Visual examples of adversarial advertisement attack generated with the COCO dataset on Stable Diffusion v1.5 for AATIM and baselines. The prompt is "Clouds soar above a tall building on a sunny day."

### 2970 A.11 POTENTIAL NEGATIVE IMPACTS, LIMITATIONS AND FUTURE WORKS

2971  
2972 In this work, the three image-caption datasets are all open-released datasets, which allow researchers  
2973 to use for non-commercial research and educational purposes. These three datasets are widely used in  
2974 the research area of generative models. All baseline codes are open-accessed resources from GitHub  
2975 and licensed under the MIT License, which only requires preservation of copyright and license  
2976 notices and includes the permissions of commercial use, modification, distribution, and private use.

2977 Our work demonstrates that text-to-image generative models can be maliciously exploited to generate  
2978 unintended advertisements. Conventional T2I advertising refers to the intentional use of a text-to-  
2979 image diffusion model by an advertiser, where the advertiser requests the inclusion of a brand (e.g.,  
2980 McDonald’s) in the prompt, and the generated image is expected to contain the branding. In contrast,  
2981 the “adversarial advertisement” problem is how to naturally embed advertisements into generated  
2982 images when the user has no advertising intention (Vice et al., 2024). An attacker may attack the  
2983 T2I DMs and implant advertisements into generated images, even when the user’s prompt has no  
2984 information about the advertised target, in order to increase the exposure of specific product brands.  
2985 To the best of our knowledge, we are the first to introduce the problem of adversarial advertisement.  
2986 We believe our work can positively impact society by providing valuable insights for future research  
2987 on the safety of T2I DMs and highlighting the importance of addressing this issue for the broader  
2988 public. Meanwhile, the technique in our paper could be misused to embed hateful or discriminatory  
2989 elements into the T2I DM. Potential mitigation includes a post-processing filter to block any unwanted  
2990 image generation.

2991 A limitation of our AATIM framework is that our advertisement-implantation method currently relies  
2992 on an English text corpus. Extending it to multilingual or even cross-lingual text-to-image generation  
2993 remains an open problem.

2994 Extending our attack into a black-box setting is a possible future direction. A practical route that  
2995 has already been explored in model-extraction literature (Carlini et al., 2024; Tamber et al., 2025;  
2996 Zhou et al., 2024; Gu et al., 2024) is to query the target API and train a high-fidelity surrogate  
2997 whose weights approximate the black-box decision function (e.g., adaptive distillation). Once such a  
2998 surrogate is obtained, our method can be applied directly to the model.

### 2999 A.12 BACKGROUND ON RANDOMIZED SMOOTHING FOR CERTIFIED ROBUSTNESS

3000  
3001 Given a classifier  $f$ , the goal of randomized smoothing for certified robustness is constructing a  
3002 smooth classifier  $g$  from  $f$ , which assigns inputs  $x \in \mathbb{R}^d$  to classes in the set  $C$ . The function  $g(x)$  is  
3003 defined by:

$$3004 \quad g(x) = \arg \max_{c \in \mathcal{Y}} \mathbb{P}(f(x + \epsilon) = c) \quad (55)$$

$$3005 \quad \text{where } \epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

3006  
3007 The classifier  $g$  identifies the class that the base classifier  $f$  will most likely predict when the input  $x$   
3008 is slightly perturbed by noise  $\epsilon$ . Let  $p_c(x)$  denote the probability that the base classifier  $f$  assigns  
3009 input  $x$  to class  $c$ , which is expressed as:

$$3010 \quad p_c(x) = \mathbb{P}_{\epsilon \sim D}(f(x + \epsilon) = c) \quad (56)$$

3011  
3012 Without loss of generality, assume that  $p_A(x)$  and  $p_B(x)$  are the probabilities for the most probable  
3013 class  $c_A$  and the second most probable class  $c_B$ , respectively. If the probability  $\mathbb{P}(f(x + \epsilon) = c_A)$   
3014 is at least  $p_A(x)$ , which in turn is greater than or equal to  $p_B(x)$ , and both of these are greater than  
3015 the maximum probability for any other class  $c \neq c_A$ , with  $\underline{p}_A(x)$  being a lower bound and  $\overline{p}_B(x)$   
3016 an upper bound, then the classifier  $g$  will consistently output  $c_A$  for any perturbation  $\delta$  in  $\mathbb{R}^d$  where  
3017  $\|\delta\|_p \leq r_p$ . Therefore, the smooth classifier  $g$  can reliably produce the correct prediction as long as  
3018 the perturbation  $\delta$  remains within the certified  $l_p$ -norm radius  $r_p$  for  $p > 0$ .

3019  
3020 **Theorem 1.6.** (Cohen et al., 2019) Let  $f : \mathbb{R}^d \rightarrow \mathcal{Y}$  be any deterministic or random function, and let  
3021  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ . Let  $g$  be defined as in (55). Suppose  $c_A \in \mathcal{Y}$  and  $\underline{p}_A, \overline{p}_B \in [0, 1]$  satisfy:

$$3022 \quad \mathbb{P}(f(x + \epsilon) = c_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{c \neq c_A} \mathbb{P}(f(x + \epsilon) = c) \quad (57)$$

3024 Then  $g(x + \delta) = c_A$  for all  $\|\delta\|_2 < R$ , where

$$3025 \quad R = \frac{\sigma}{2} (\Phi^{-1}(p_A) - \Phi^{-1}(p_B)) \quad (58)$$

3026  $\Phi^{-1}$  is the inverse of the standard Gaussian CDF. Please refer to the original paper Cohen et al.  
3027 (2019) for detailed proof.

3028 Recent works (Kumar et al., 2020; Yang et al., 2020; Mohapatra et al., 2020) have revealed that  
3029 the largest certified radius  $r_p$  for randomized smoothing against  $l_p$ -norm adversarial threats scales  
3030 inversely with  $d^{\frac{1}{2} - \frac{1}{p}}$ , where  $d$  denotes the input dimension. Specifically, for a Gaussian distribution  
3031 with variance  $\sigma^2$ , the upper bound of  $r_p$  is given by (Kumar et al., 2020):

$$3032 \quad r_p = \frac{\sigma}{2d^{\frac{1}{2} - \frac{1}{p}}} (\Phi^{-1}(p_A(x)) - \Phi^{-1}(p_B(x))) \quad (59)$$

3033 In this context,  $\sigma$  functions as a hyperparameter to balance robustness and accuracy within the model  
3034  $g$ . It's noted that as the dimension  $d$  increases, particularly when  $p > 2$ , the upper bound of  $r_p$   
3035 significantly decreases, rendering the certified radius extremely small for high-dimensional spaces.  
3036 Consequently, this weakens robustness against  $l_p$ -norm adversarial attacks in high-dimensional  
3037 contexts.

### 3038 A.13 THE SOLUTION OF MULTIVARIATE CONTINUOUS SCALED PHASE-TYPE WITH LÉVY 3039 DISTRIBUTION

3040 The partial derivatives with respect to the parameters are computed below.

$$3041 \quad \frac{\partial L}{\partial \alpha} = \frac{P_A(x)e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathbf{A} \mathbf{1}}{-1 + \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathbf{A} \mathbf{1}} + \frac{1 - P_A(x)}{\alpha} = 0, \quad (60)$$

$$3042 \quad \frac{\partial L}{\partial \mathbf{B}} = \frac{P_A(x)\alpha e^{\eta \mathbf{B} \sqrt{x}} \eta \sqrt{x} \mathbf{D} \mathbf{A} \mathbf{1}}{-1 + \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathbf{A} \mathbf{1}} + \eta \sqrt{x} (1 - P_A(x)) = 0, \quad (61)$$

$$3043 \quad \frac{\partial L}{\partial \mathbf{D}} = \frac{P_A(x)\alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{A} \mathbf{1}}{-1 + \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{A} \mathbf{1}} + \frac{(1 - P_A(x))\alpha e^{\eta \mathbf{B} x} \mathbf{A} \mathbf{1}}{\alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{A} \mathbf{1}} = 0, \quad (62)$$

$$3044 \quad \frac{\partial L}{\partial \eta} = \frac{P_A(x)\alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{B} \sqrt{x} \mathbf{D} \mathbf{A} \mathbf{1}}{-1 + \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathbf{A} \mathbf{1}} + \mathbf{B} \sqrt{x} (1 - P_A(x)) = 0, \quad (63)$$

$$3045 \quad \frac{\partial L}{\partial \mathbf{A}} = \frac{P_A(x)\alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathbf{1}}{-1 + \alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathbf{1}} + \frac{(1 - P_A(x))\alpha e^{\eta \mathbf{B} x} \mathbf{D} \mathbf{1}}{\alpha e^{\eta \mathbf{B} \sqrt{x}} \mathbf{D} \mathbf{1}} = 0. \quad (64)$$

3046 The solution to the above equations are

$$3047 \quad \alpha = \mathbf{1}^{-1} \mathbf{A}^{-1} \mathbf{D}^{-1} e^{-\eta \mathbf{B} \sqrt{x}} (1 - P_A(x)), \quad (65)$$

$$3048 \quad \mathbf{B} = \frac{\log(\alpha^{-1} (1 - P_A(x)) \mathbf{1}^{-1} \mathbf{A}^{-1} \mathbf{D}^{-1})}{\eta \sqrt{x}}, \quad (66)$$

$$3049 \quad \mathbf{D} = e^{-\eta \mathbf{B} \sqrt{x}} \alpha^{-1} (1 - P_A(x)) \mathbf{1}^{-1} \mathbf{A}^{-1}, \quad (67)$$

$$3050 \quad \eta = \frac{\log(\alpha^{-1} (1 - P_A(x)) \mathbf{1}^{-1} \mathbf{A}^{-1} \mathbf{D}^{-1})}{\sqrt{x} \mathbf{B}}, \quad (68)$$

$$3051 \quad \mathbf{A} = \mathbf{D}^{-1} e^{-\eta \mathbf{B} \sqrt{x}} \alpha^{-1} (1 - P_A(x)) \mathbf{1}^{-1}, \quad (69)$$

3052 where the inverse notation is used to represent vectors  $\alpha^{-1}$  and  $\mathbf{1}^{-1}$  such that  $\mathbf{1}^{-1} \times \mathbf{1} = 1$  and  
3053  $\alpha \times \alpha^{-1} = 1$ .

### 3054 A.14 THE USE OF LARGE LANGUAGE MODELS

3055 In this submission, we used an LLM solely to polish the writing and correct grammatical errors.

### 3078 A.15 COMPUTATIONAL COST AND OFFLINE NATURE

3079  
3080  
3081  
3082 First, the MCPHL module introduces only a small number of parameters. Let  $m$  denote the number  
3083 of MCPHL states. The sub-intensity matrix  $B \in \mathbb{R}^{m \times m}$ , together with the diagonal matrices  $D$   
3084 and  $A$  and the initial vector  $\alpha$ , yields  $O(m^2)$  parameters in total, whereas the text encoder typically  
3085 contains hundreds of millions of parameters. Hence, the additional parameter footprint of MCPHL is  
3086 quadratic in the small constant  $m$  and negligible compared to the backbone encoder.

3087 Second, AATIM is an offline attack: the adversary optimizes the adversarial advertisement and the  
3088 MCPHL density prior to deployment, without any real-time latency constraints, and only needs to  
3089 upload the attacked checkpoints. As a result, computational efficiency does not pose a practical  
3090 limitation on the applicability of our method.

3091 Third, in realistic advertising scenarios it is reasonable to assume that each advertiser only needs  
3092 to advertise a single brand (Romano, 2005; Beard, 2013; Jagpal & Jagpal, 2008; Bai et al., 2021).  
3093 Under this assumption, fitting the MCPHL distribution is required only once per brand: the learned  
3094 density and the corresponding advertisement tokens can be reused for all future prompts related to  
3095 that brand, without re-estimating the distribution. As a result, the one-time cost of fitting MCPHL  
3096 for an attacker is trivial in terms of computational efficiency. Empirically, the total time used for an  
3097 advertisement injection on SD1.5 is around 32 minutes for a single NVIDIA H100 GPU.

### 3098 3099 3100 3101 3102 3103 A.16 OVERALL THREE-STAGE FORMULATION OF AATIM

3104  
3105  
3106  
3107 Conceptually, our method consists of three stages: fitting the MCPHL, attacking the encoder,  
3108 and performing masked mollification. In the first stage, we estimate an MCPHL density  $p()$  on  
3109 sentence embeddings of real advertisement prompts. **Definition 3.1** introduces the continuous scaled  
3110 phase-type family as the basic parametric form we build on; **Definition 3.2** introduces the Lévy  
3111 distribution, which we use as the positive scaling variable to capture the heavy-tail nature; **Definition**  
3112 **3.3** combines these two pieces into the one-dimensional continuous scaled phase-type with Lévy  
3113 (CPHL) distribution; and **Definition 3.4** extends CPHL to the multivariate MCPHL distribution  
3114 used to model the joint distribution of prompt embeddings. **Theorem 3.5** then shows that, under  
3115 our estimation objective, the resulting MCPHL CDF  $Q_A(x)$  converges to the empirical distribution  
3116  $P_A(x)$  of real advertisement embeddings. In summary, Definitions 3.1–3.4 provide a step-by-step  
3117 theoretical construction of the MCPHL, while **Eq. (6)** gives the objective function used to fit the  
3118 MCPHL in our implementation.

3118  
3119 In the second stage, given the trained MCPHL and its pdf  $p$ , we optimize the encoder using the  
3120 objective in **Eq. (7)**. This is the actual loss used to train the attack encoder in our implementation.  
3121 The pdf  $p$  learned in Stage 1 is used in the second term of **Eq. (7)** as a density regularizer. The  
3122 MCPHL parameters are fixed in this stage and are not updated jointly with the encoder.

3123  
3124 In the third stage, after we obtain the attacked encoder (denoted by  $f(w)$ ) from Stage 2, we apply the  
3125 masked mollification procedure in Section 4 as a post-processing step to obtain a smoothed encoder  
3126  $g(w)$ . **Definition 4.1** specifies the main mollification step, where the original function is convolved  
3127 with a Friedrichs kernel. The subsequent definitions and theorems provide the theoretical framework  
3128 that yields dimension-invariant robustness guarantees for our masked mollification. In particular,  
3129 **Definition 4.2** introduces the Hadamard directional derivative, which is then used in **Theorem 4.3** to  
3130 show that the  $\ell_p$  norm is Hadamard-directionally differentiable. Given this differentiability, we derive  
3131 the Lipschitz constant of the mollified function  $G$  in **Theorem 4.4**, which then allows us to obtain a  
dimension-independent certified radius in **Theorem 4.5**. In summary, Definition 4.1 describes the  
mollification step used in our implementation, while Definitions 4.2 and Theorems 4.3–4.5 provide  
the theoretical backbone that explains and justifies our technical contribution.

If we try to bring these three related but not jointly optimized stages together in a unified formulation, it can be written as follows:

$$\begin{aligned}
 \mathcal{L}_{\text{all}}(\theta, w) = & \underbrace{\mathbb{E}_{x \sim \mathcal{D}} [P_A(x) \log Q_{A, \theta}(x) + (1 - P_A(x)) \log (1 - Q_{A, \theta}(x))]}_{\text{(i) fitting the MCPHL (Eq. 6)}} \\
 & + \lambda_A \underbrace{\mathbb{E}_{(s, \hat{s}) \sim \mathcal{D}_A} \left\| E_{G(w)}(s) - E_{f, G(w)}(\hat{s}) \right\|_2^2}_{\text{(ii) encoder attack loss (cf. Eq. 7)}} \\
 & - \lambda_M \underbrace{\mathbb{E}_{s \sim \mathcal{D}_A} \log p(E_{G(w)}(s))}_{\text{(iii) MCPHL density regularizer, using } p}
 \end{aligned} \tag{70}$$

where  $\theta = (\alpha, \eta, \mathbf{B}, \mathbf{D}, \mathcal{A})$  denotes the MCPHL parameters,  $\lambda_A, \lambda_M$  are hyperparameters weighting the attack and density terms,  $p$  is its pdf, and the mollified encoder weights  $G(w)$  are given by

$$G(w) = \int F(w - \text{Mask}(w) \odot u) \varphi_\sigma(u) du \quad (\text{cf. Eq. 8}). \tag{71}$$

However, as we have emphasized above, these three stages are not trained jointly in practice, so forcing them into a single objective function would hurt readability and could be misleading.

#### A.17 CERTIFIED ROBUSTNESS IN PARAMETER SPACE

There are already many works that explicitly applies randomized smoothing or related certification techniques to neural network parameter spaces instead of input space. For example: Bansal et al. (2022) apply randomized smoothing directly in the model parameter space and derive certified  $l_2$ -radii against weight perturbations; Gan et al. (2023) analyze vulnerability in the parameter space and optimize watermark robustness under parametric changes; Weng et al. (2020) develop formal robustness guarantees under adversarial perturbations of network weights; Fischer et al. (2020) extend randomized smoothing to certify robustness with respect to transformation parameters rather than input space.

From the above lines of work, we can find substantial evidence to justify that small parameter changes during fine-tuning can be treated analogously to adversarial input perturbations. Weng et al. (2020) explicitly defines a threat model in which the parameters are adversarially perturbed within a norm ball and derives certified bounds: if  $\|\Delta w\| \leq \varepsilon$ , then the prediction cannot change. This is structurally the same as certification in input space, but with weights instead of the input  $x$ . Tsai et al. (2021) further formalizes generalization and adversarial robustness of neural networks to weight perturbations, providing margin and generalization bounds under norm-bounded weight noise and proposing a theory-driven objective that improves robustness to such perturbations. Savva et al. (2023) discusses robustness bounds to weight perturbations and empirically shows that misclassifications can be triggered by small changes in the parameters, very similar to adversarial input perturbations.

Moreover, all previous works (Bansal et al., 2022; Gan et al., 2023; Weng et al., 2020; Fischer et al., 2020; Tsai et al., 2021; Savva et al., 2023) study deep, highly nonlinear neural networks and derive robustness margins, certified bounds, or generalization guarantees under norm-bounded weight perturbations, rather than for simple linear models. This line of work shows that certified robustness analysis does not rely on linearity of the model, and further justifies our use of the analogy between perturbations in parameter space and perturbations in input space.

Compared with prior work, our contribution does not lie in drawing an analogy between parameter perturbations and input-space perturbations, since that has already been justified a lot. Instead, our contribution is the introduction of a masked, dimension-invariant mollification scheme together with the derivation of a dimension-invariant certified radius, which makes this approach better suited for text-to-image models with large number of parameters.

#### A.18 GENERATING HIGH-QUALITY BRANDED DATA WITH LLM

Our method does not require any preexisting large corpus of branded sentences. Instead, given a random MS-COCO caption and a brand name (e.g., ‘‘McDonald’s’’), we query a general-purpose LLM (e.g., ChatGPT) to generate a single fluent sentence that naturally incorporates the brand.

3186 Modern LLMs are explicitly designed to produce high-quality text with lexical constraints, many  
3187 existing works in top ML/NLP venues treats LLMs as reliable generators for synthetic labeled  
3188 text data (Ouyang et al., 2022; Yu et al., 2024; Ma et al., 2024). This suggest that generating  
3189 branded sentence with LLM is a standard and well-validated practice, and our approach relies on  
3190 this generation rather than some large, high quality branded text corpus. In addition, the MCPHL  
3191 density in our method is learned from data with a formal convergence guarantee: Theorem 3.5 shows  
3192 that, under the estimation objective in Eq. (6), the learned MCPHL CDF  $Q_A(x)$  converges to the  
3193 empirical distribution  $P_A(x)$  of advertisement embeddings. This provides a principled justification  
3194 that our density estimation are accurate on branded prompts.

3195  
3196  
3197  
3198  
3199  
3200  
3201  
3202  
3203  
3204  
3205  
3206  
3207  
3208  
3209  
3210  
3211  
3212  
3213  
3214  
3215  
3216  
3217  
3218  
3219  
3220  
3221  
3222  
3223  
3224  
3225  
3226  
3227  
3228  
3229  
3230  
3231  
3232  
3233  
3234  
3235  
3236  
3237  
3238  
3239