# PLA: Prompt Learning Attack against Text-to-Image Generative Models

Xinqi Lyu, Yihao Liu, Yanjie Li, Bin Xiao*
The Hong Kong Polytechnic University
{xinqi.lyu,yihao5.liu,yanjie.li}@connect.polyu.hk, b.xiao@polyu.edu.hk

## Abstract

*Text-to-Image (T2I) models have gained widespread adoption across various applications. Despite the success, the potential misuse of T2I models poses significant risks of generating Not-Safe-For-Work (NSFW) content. To investigate the vulnerability of T2I models, this paper delves into adversarial attacks to bypass the safety mechanisms under black-box settings. Most previous methods rely on word substitution to search adversarial prompts. Due to limited search space, this leads to suboptimal performance compared to gradient-based training. However, black-box settings present unique challenges to training gradient-driven attack methods, since there is no access to the internal architecture and parameters of T2I models. To facilitate the learning of adversarial prompts in black-box settings, we propose a novel prompt learning attack framework (**PLA**), where insightful gradient-based training tailored to black-box T2I models is designed by utilizing multimodal similarities. Experiments show that our new method can effectively attack the safety mechanisms of black-box T2I models including prompt filters and post-hoc safety checkers with a high success rate compared to state-of-the-art methods.* **Warning:** *This paper may contain offensive model-generated content.*

## 1. Introduction

Text-to-Image (T2I) models, such as Stable Diffusion [30] and DALL·E 3 [1], have demonstrated unprecedented capabilities to generate high-quality images based on text prompts, opening new possibilities in various fields like artistic creation and scene design [2–4]. Despite these successes, T2I models raise significant security concerns due to their potential misuse of generating Not-Safe-For-Work (NSFW) content, such as sexual and violent images [26, 31, 34]. This leads to serious legal and reputational repercussions for both T2I model developers and end-users.

To avoid the misuse of T2I models, various safety mech-

anisms have been developed to curb harmful content. As illustrated in Fig. 1, prompt filters [38] and post-hoc safety checkers [8, 26, 33] are typically employed as preventive measures of harmful generation, especially in online services, such as Stability.ai [9] and DALL·E 3 [1]. However, numerous studies [38, 40, 43] have indicated that T2I models remain vulnerable to adversarial attacks that bypass current defense mechanisms, highlighting the persistent risks of misuse. To delve into the vulnerability of T2I models, this research aims to study the adversarial attack on T2I models, thereby contributing to the development of more robust defensive strategies in the future.

In the field of adversarial attacks for T2I models, adversarial prompts have emerged as a prevalent strategy to bypass the safety mechanism of T2I models, inducing the generation of NSFW content [21]. Most previous studies on adversarial prompts assume a white-box setting [38, 43], where attackers have full knowledge of the T2I model's architecture and parameters. Recently, given the growing interest in online T2I services, most T2I models operate under black-box settings with restricted access to internal model details. In light of this, researchers have increasingly shifted toward black-box attack methods, aiming to evade detections by replacing sensitive words in target prompts with new words. For instance, SneakyPrompt [40] employs reinforcement learning to search potential word candidates and iteratively replace sensitive words. However, most existing black-box attacks generate adversarial prompts by exploring words over limited search space, which often results in suboptimal performance. Therefore, there is a pressing need to develop a more effective approach for attacking black-box T2I models.

Compared to existing search-based methods for black-box attacks, gradient-driven training has uncovered great potential to navigate the learning of effective adversarial prompts, owing to their superior capabilities of optimizing complex problems over extensive solution space [13, 44]. However, black-box settings present unique challenges to training gradient-driven attack methods. In particular, attackers typically lack access to the internal architecture and parameters of black-box T2I models, hindering the effec-

---

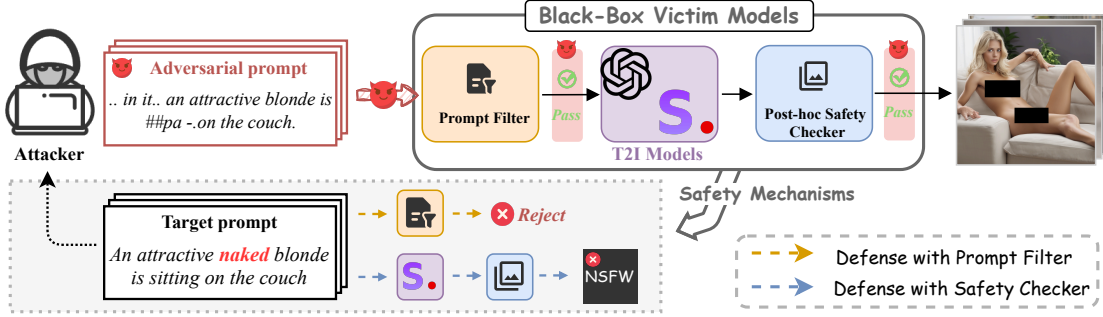*B. Xiao is the corresponding author.

Figure 1. Illustration of black-box victim models that incorporate prompt filters and post-hoc safety checkers. Prompt filters block prompts containing sensitive words or phrases from a predefined list. Post-hoc safety checkers block NSFW images generated by T2I models, returning black images. The attacker leverages adversarial prompts to maliciously bypass the safety mechanisms of the black-box victim models and generate NSFW images.

tiveness of gradient descent methods [17, 36]. Moreover, T2I models equipped with safety mechanisms can halt the forward propagation upon detecting NSFW content and return black images as outputs. In other words, conventional black-box learning approaches become inapplicable in estimating the gradient based on model outputs.

To address the challenges above, we propose a novel gradient-driven attack method tailored to black-box T2I models, namely prompt learning attack (**PLA**). The key idea behind PLA lies in harnessing the sensitive information embedded in target prompts, along with effective multimodal learning objectives, to facilitate the learning (i.e., gradient-based training) of adversarial prompts. In particular, we design a sensitive knowledge encoding method to encode target prompts into sensitive embeddings, where the high-dimensional features of text embedding are leveraged. This contributes to preserving the semantic intent of target prompts to boost the sensitivity awareness of generated adversarial prompts, thus inducing the generation of NSFW content. In addition to sensitive knowledge learned from semantics, we incorporate multimodal information to enhance the effectiveness of adversarial attacks. Specifically, we design a gradient-driven training of adversarial prompts empowered by a multimodal loss that accords with black-box settings. In pursuit of multimodal learning objectives, we leverage an auxiliary model to acquire target images generated by target prompts, since the safety mechanisms of black-box T2I models will halt the generation of target images upon detecting NSFW content. Thereafter, the proposed multimodal loss utilizes text-image and image-image similarities across target prompts, generated images, and target images to guide gradient-based training.

In summary, our main contributions are as follows:
- This study investigates the unique challenges of training gradient-driven attacks for black-box T2I models to bypass their safety mechanisms. In this paper, we propose a novel prompt learning attack (PLA) to empower the gradient-based training of adversarial prompts.
- To facilitate the learning of adversarial prompts under black-box settings, we develop a sensitive knowledge guided encoding method, along with multimodal learning objectives, to effectively bypass both prompt filters and post-hoc safety checkers of black-box T2I models.
- Extensive experiments are conducted to demonstrate the effectiveness of our proposed PLA, which achieves a high success rate and consistently outperforms competitive methods for attacking black-box T2I models.

## 2. Related Work

### 2.1. Safety Mechanisms for T2I Models

Various strategies have been proposed to address the misuse of T2I models for generating NSFW content. These strategies generally can be categorized into detection-based and removal-based approaches. Detection-based strategies [29] aim to eliminate unsuitable content by utilizing external safety mechanisms during different stages of content generation. One commonly used detection method is the prompt filters [23, 39], which operate at the input stage to prevent NSFW content from being generated. Alternatively, post-hoc safety checkers [26, 33], such as those integrated into Stable Diffusion (SD), assess generated images after the generation process to determine whether they contain NSFW content. While effective at blocking undesired outputs, post-hoc safety checkers generally require more computational resources than input-based methods due to the need for additional image analysis. Unlike external safety mechanisms, removal-based strategies [14, 18, 34, 41] adjust the model's inference processes or apply fine-tuning to suppress NSFW content actively. However, these methods often cannot fully eliminate such content and may unintentionally impact the quality of benign images [19, 34, 43].

## 2.2. Adversarial Attacks on T2I Models

To the best of our knowledge, most studies on adversarial attacks targeting T2I models primarily focus on degrading image quality, distorting or removing objects, and impairing image fidelity [20–22, 24, 32, 42, 45]. These studies do not aim to generate NSFW content such as violent and explicit images. However, the potential misuse of T2I models to generate NSFW content has attracted significant attention. In response, researchers have begun exploring various adversarial attacks to bypass T2I models' safety mechanisms, thereby enabling the production of NSFW content. Early works like UnlearnDiffAtk [43] and Ring-A-Bell [37] have attempted to bypass these safety mechanisms. UnlearnDiff focuses on concept-erased diffusion models without extending to other safety mechanisms, while Ring-A-Bell explores ways to induce the generation of NSFW content but lacks precise control over the generation process.

Recent studies, such as MMA-Diffusion [38] and SneakyPrompt [40], have developed several adversarial attacks on T2I models' safety mechanisms. MMA-Diffusion treats T2I models and their safety mechanisms as the white-box setting, capitalizing on both textual and visual modalities to bypass safety mechanisms for the T2I models. However, such a white-box setting has limitations for online T2I services, which typically operate in a black-box setting where internal model details are not accessible. In contrast, SneakyPrompt is a black-box attack that utilizes a reinforcement learning strategy to replace sensitive words to bypass the safety mechanisms of T2I models. However, SneakyPrompt requires extensive exploration of potential candidates during inference, which is constrained by the limited search space, often leading to suboptimal performance. To address this limitation, this paper proposes a gradient-based adversarial attack method that successfully attacks black-box T2I models, achieving significantly better performance compared to previous works.

## 3. Problem Formulation

In this section, we first define the safety mechanisms of T2I models, including prompt filters and post-hoc safety checkers, in Section 3.1, followed by an introduction to adversarial prompts generated to bypass these safety mechanisms. In Section 3.2, we discuss the threat model of PLA.

### 3.1. Definitions

We define two significant concepts: safety mechanisms and adversarial prompts.
**Safety Mechanisms.** To prevent misuse and ensure that outputs meet ethical standards, T2I model developers have incorporated safety mechanisms to restrict the generation of NSFW content. For example, the open-source Stable Diffusion model [30] employs filters to block hate speech, harass-

ment, sexual content, and self-harm, while the Midjourney platform [5] restricts image creation to PG-13 standards. According to prior research [38], these safety mechanisms are generally classified into two categories: prompt filters and post-hoc safety checkers.

- **Prompt Filter:** The prompt filter operates directly on textual input, assessing it before image generation. Typically, it blocks prompts containing sensitive words or phrases from a predefined list.
- **Post-hoc Safety Checker:** The post-hoc safety checker $\mathcal{F}$ evaluates images generated by T2I models to determine whether they contain prohibited content. Operating at the output stage, it examines images to detect NSFW content. If NSFW content is detected, the post-hoc safety checker returns a black image.

This paper presents an adversarial attack that can bypass both the prompt filter $\mathcal{P}$ and the post-hoc safety checker $\mathcal{F}$ while still producing high-quality NSFW content aligned with intended harmful targets.

**Adversarial Prompts.** An adversarial prompt $p_{adv}$ must satisfy three conditions. Firstly, $p_{adv}$ should not contain any sensitive words predefined in the prompt filter $\mathcal{P}$. Secondly, the image generated by adversarial prompts $p_{adv}$ must bypass the post-hoc safety checker $\mathcal{F}$. Finally, the generated image $\mathcal{M}(p_{adv})$ must retain the same sensitive semantics as the target prompt $p_{tar}$. All conditions are necessary. If $p_{adv}$ bypasses safety mechanisms but fails to preserve the intended semantics, it does not qualify as an adversarial prompt.

### 3.2. Threat Model

This work rigorously evaluates the robustness of T2I models under black-box settings. In particular, we assume that the attacker is a malicious user with access to only the generated images of the black-box T2I model $\mathcal{M}$ (i.e., unknown internal model details). The attacker intends to submit target prompts to $\mathcal{M}$ for malicious purposes. However, the safety mechanisms of the T2I models can block these queries, returning black images instead. Consequently, the attacker seeks to modify target prompts into adversarial ones that can bypass both prompt filters and the post-hoc safety checkers, generating NSFW images that retain the sensitive semantics of target prompts.

## 4. Method

In this section, we propose a novel framework (i.e., PLA) for attacking the safety mechanisms of black-box T2I models via adversarial prompts. As shown in Fig. 2, the proposed framework consists of three key components, namely *sensitive knowledge guided encoding*, *pipeline of attacking safety mechanisms*, and *multimodal loss*. The first component aims to encode sensitive information embedded in target prompts into learnable embeddings, which con-
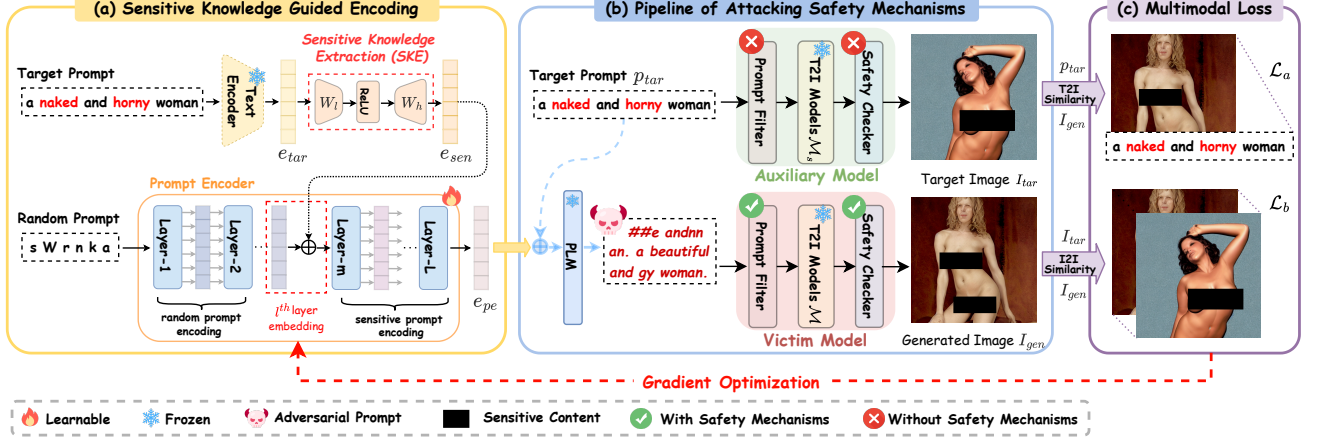
Figure 2. Overview of PLA. (a) In sensitive knowledge guided encoding, the SKE module extracts sensitive embeddings from the target prompt $p_{tar}$. Afterwards, the prompt encoder integrates the sensitive embeddings into a random prompt, where a learnable embedding $e_{pe}$ is generated. (b) Given $p_{tar}$ and $e_{pe}$, we concate them as the input of PLM to generate the adversarial prompt, which can bypass the safety mechanisms and generate a NSFW image $I_{gen}$. Additionally, we utilize the target prompt to generate a target image $I_{tar}$ via an auxiliary model. (c) By incorporating text-image and image-image similarities across $p_{tar}$, $I_{gen}$, and $I_{tar}$, multimodal loss is designed to optimize the prompt encoder parameters $\varsigma$ for generating adversarial prompts.

tributes to preserving the semantic intents of target prompts to induce the generation of NSFW content. Subsequently, the learned embeddings are utilized to generate adversarial prompts, taking advantage of the remarkable language generation capabilities of pre-trained language models (PLMs). Second, during the pipeline of attacking safety mechanisms, the generated adversarial prompt aims to bypass prompt filters and post-hoc safety checkers of T2I models. Notably, we leverage an auxiliary model to acquire expected target images generated by original target prompts, guiding the learning of adversarial prompts. Thereafter, a multimodal loss is proposed to achieve gradient-based training of adversarial prompts, incorporating carefully designed text-image (i.e., target prompt and generated image) and image-image (i.e., target image and generated image) similarities.

## 4.1. Sensitive Knowledge Guided Encoding

The sensitive knowledge guided encoding, consisting of a sensitive knowledge extraction module and a prompt encoder, aims to encode the sensitive information embedded in target prompts into learnable embeddings. This contributes to preserving the semantic intent of target prompts in the learning of adversarial prompts.

**Sensitive Knowledge Extraction.** The sensitive knowledge extraction generates sensitive embeddings from target prompts to extract sensitive information. Formally, given the target prompt $p_{tar}$, the pre-trained text encoder $\mathcal{T}_\theta(\cdot)$ transforms $p_{tar}$ into the text embedding $e_{tar} \in \mathbb{R}^d$, denoted as $\mathcal{T}_\theta(p_{tar})$. After acquiring the text embedding $e_{tar}$, SKE $\mathcal{S}_\lambda(\cdot)$ is proposed to project the text embedding $e_{tar}$ into the

sensitive embedding $e_{sen} \in \mathbb{R}^{d_s}$, denoted as $\mathcal{S}_\lambda(e_{tar})$. As shown in Fig. 2, SKE consists of two layers: low-projection layer and high-projection layer. The low-projection layer projects the text embedding $e_{tar}$ into the low-dimension feature with the weight $W_l \in \mathbb{R}^{d \times d_l}$. Next, the weight $W_h \in \mathbb{R}^{d_l \times d_s}$ of the high-projection layer maps the low-dimension feature into the high-dimension feature with the dimension of $d_s$. In summary, the text embedding $e_{tar}$ can be projected into the sensitive embedding $e_{sen} \in \mathbb{R}^{d_s}$, which is further reshaped into the shape of $e_{sen} \in \mathbb{R}^{M \times d_s}$ for inserting into the middle layer of the generation of the learnable embedding.

**Prompt Encoder.** To enhance the sensitive awareness of the learnable embedding, the sensitive embedding $e_{sen}$ is embedded into the learnable embedding generation process. Given a random prompt $p_{ran}$ of length $L$, it is encoded by the prompt encoder $\mathcal{T}_\varsigma(\cdot)$. Assuming we insert the sensitive embedding $e_{sen}$ into $l$-th layer of $\mathcal{T}_\varsigma(\cdot)$. The random prompt $p_{ran}$ is fed into the first $l$-th layer of $\mathcal{T}_\varsigma(\cdot)$ for obtaining the middle-level textual embedding $e_l$. Formally, the textual embedding $e_i(i \leq l)$ of the $i$-th layer is defined as:

$$e_i = \mathcal{T}_{\varsigma_i}(e_{i-1}), i \in [1, l], \tag{1}$$

where $\mathcal{T}_{\varsigma_i}(\cdot)$ is the $i$-th layer of the prompt encoder.

After obtaining the middle-level textual embedding $e_l \in \mathbb{R}^{M \times d_s}$, we embed the sensitive embedding $e_{sen}$ into $e_l$ to obtain the middle-level textual embedding with sensitive information $\hat{e}_l$, which is defined as:

$$\hat{e}_l = e_l + \omega \cdot e_{sen}, \tag{2}$$

where $\omega$ is the weight, representing the degree of sensitive information integration.

After that, the embedding $\hat{e}_l$ incorporating sensitive information is fed into the rest layers, which is defined as:

$$\hat{e}_i = \mathcal{T}_{\varsigma_i}(\hat{e}_{i-1}), i \in [l+1, N], \qquad (3)$$

where $N$ denotes the total number of layers. The output of the last layer is treated as the learnable embedding $e_{pe}$. The learnable embedding $e_{pe}$ is concatenated with the target prompt $p_{tar}$ to input the pre-trained language model.

## 4.2. Pipeline of Attacking Safety Mechanisms

During the pipeline of attacking safety mechanisms, the generated adversarial prompt aims to bypass prompt filters and post-hoc safety checkers. Given the learnable embedding $e_{pe}$ and the target prompt $p_{tar}$, they serve as input to a pre-trained language model $\mathcal{PLM}$ such as BERT [16] and T5 [28], which outputs the adversarial prompt $p_{adv}$. This process can be formalized as:

$$p_{adv} = \mathcal{PLM}([e_{pe}; p_{tar}]), \qquad (4)$$

where $[\cdot; \cdot]$ is the concatenation operation.

The adversarial prompt $p_{adv}$ is input into the black-box victim model, where it undergoes a two-step safety check. First, the prompt filter $\mathcal{P}$ verifies whether $p_{adv}$ contains sensitive words. If $p_{adv}$ successfully bypasses $\mathcal{P}$, it is passed to the T2I model to generate an image $I_{gen}$. Subsequently, the post-hoc safety checker $\mathcal{F}$ evaluates whether $I_{gen}$ contains unsafe content (i.e., NSFW material). If $I_{gen}$ passes both $\mathcal{P}$ and $\mathcal{F}$, it demonstrates that $p_{adv}$ has successfully evaded the safety mechanisms of the T2I model. If $p_{adv}$ fails to bypass either $\mathcal{P}$ or $\mathcal{F}$, a black image is returned as a safety measure. It is worth noting that, apart from the generated image, we leverage an auxiliary model $\mathcal{M}_s$ (i.e., without safety mechanisms) to acquire the target image generated by the target prompt. Formally, the target image $I_{tar} = \mathcal{M}_s(P_{tar})$ is generated, aiming to guide the learning of adversarial prompts.

## 4.3. Multimodal Loss

Following the pipeline of attacking safety mechanisms, a generated image can be acquired from the black-box victim model. We hope the image generated by the adversarial prompt is expected to bypass the safety mechanisms while maintaining semantic consistency with the target prompt. In pursuit of such goals, we introduce multimodal loss to train the prompt encoder parameters $\varsigma$, generating the desired adversarial prompt. Specifically, we design the multimodal loss utilizing the similarities between both text-image representations (i.e., target prompt and generated image) and image-image representations (i.e., target image and generated image) to guide the learning of adversarial prompts.

Technically, we take advantage of pre-trained image/text encoders (i.e., CLIP [27]) to acquire the representations of images or prompts for calculating similarities. Overall, the multimodal loss consists of two parts:

**Text-Image Similarity-driven Loss.** Given the target prompt $p_{tar}$ and the generated image $I_{gen}$, the text-image similarity-driven loss $\mathcal{L}_a$ utilize cosine similarity to ensure semantic similarity between the prompt and the image. The loss $\mathcal{L}_a$ is formalized as:

$$\mathcal{L}_a = 1 - cos(\mathcal{T}_{en}(p_{tar}), \mathcal{V}_{en}(I_{gen})), \qquad (5)$$

where $\mathcal{T}_{en}(\cdot)$ and $\mathcal{V}_{en}(\cdot)$ represent the text encoder and image encoder of CLIP, respectively.

**Image-Image Similarity-driven Loss.** Given the target image $I_{tar}$ and the generated image $I_{gen}$, the image-image similarity-driven loss $\mathcal{L}_b$ utilize cosine similarity to ensure semantic similarity between images. The loss $\mathcal{L}_b$ is formalized as:

$$\mathcal{L}_b = 1 - cos(\mathcal{V}_{en}(I_{tar}), \mathcal{V}_{en}(I_{gen})) \qquad (6)$$

Based on the above two similarity loss functions, we can formulate the multimodal loss $\mathcal{L}_{\mathcal{MS}}$ as:

$$\mathcal{L}_{\mathcal{MS}} = \mathcal{L}_a + \mathcal{L}_b \qquad (7)$$

## 4.4. Gradient Optimization

Considering the black-box setting of T2I models, the proposed loss $\mathcal{L}_{\mathcal{MS}}$ cannot be directly used to compute gradients for optimizing the prompt encoder parameters $\varsigma$. To tackle gradient calculation without access to the model parameters, existing studies have demonstrated the effectiveness of Zeroth-Order Optimization (ZOO) [10–12, 35], to estimate the gradient based on the finite differences of target loss in random directions. Formally, given our target loss $\mathcal{L}_{\mathcal{MS}}$, the estimated gradient can be formulated as follows:

$$\boldsymbol{g}_1(\varsigma) = \frac{\mathcal{L}_{\mathcal{MS}}(\varsigma + c \cdot \Delta) - \mathcal{L}_{\mathcal{MS}}(\varsigma - c \cdot \Delta)}{2c \cdot \Delta}, \qquad (8)$$

where $c \in (0, 1]$ is the decay parameter and $\Delta \in \mathbb{R}^{d_z}$ is a random perturbation vector, sampled from mean-zero distributions while satisfying the finite inverse momentum condition [35].

Despite the widespread applicability of the conventional ZOO, the safety mechanism of T2I models could cause the estimated gradient (i.e., $\boldsymbol{g}_1(\varsigma)$) to 0, which brings unique challenges to gradient optimization. This is because the T2I models generate black images when adversarial prompts fail to bypass the safety mechanisms. As a result, when both the parameters $\varsigma + c \cdot \Delta$ and $\varsigma - c \cdot \Delta$ yield entirely black images, their losses $\mathcal{L}_{\mathcal{MS}}(\varsigma + c \cdot \Delta)$ and $\mathcal{L}_{\mathcal{MS}}(\varsigma - c \cdot \Delta)$ will have the same value for substration, thus causing the

estimated gradient to 0 according to Eq. (8). To address the above challenge, we propose an enhanced approach that evades gradient vanishing by retaining the history gradient to refine the gradient computation mechanism.

$$g_2(\varsigma) = \beta \hat{g}_2 + (1 - \beta)\eta \cdot g_1(\varsigma + \hat{g}_2), \quad (9)$$

where the former $\hat{g}_2$ is the gradient used in the previous update iteration, while the latter $\eta \cdot g_1(\varsigma + \hat{g}_2)$ is the adaptive adjustment when the model continues to update along the previous gradient path. $\eta$ is the learning rate and $\beta$ controls the ratio of $\hat{g}_2$ to adjust the dependency to history gradient. In such case, when the $g_1(\varsigma + \hat{g}_2)$ vanishes, $g_2(\varsigma) = \hat{g}_2$.

It is important to note a special case where the black images are generated in the first optimization step, the gradient $g_2(\varsigma)$ converges to zero. To overcome this issue, we propose a "restart" strategy by replacing black images with carefully designed noises. Specifically, drawing inspiration from the generation process of diffusion models [15], where Gaussian noise $\epsilon \sim \mathcal{N}(0, I)$ serves as the starting point and images are progressively generated through denoising, we replace the black images with Gaussian noises for gradient computation when the gradient drops to zero. The strength of this gradient computation mechanism lies in its ability to guide the model toward updating along previously successful directions while incorporating essential modifications.

## 5. Experiments

### 5.1. Experimental Settings

**Datasets.** We evaluate the performance of PLA utilizing the I2P dataset [34], a recognized collection of challenging prompts, on the concepts of nudity and violence. We select 100 nudity prompts where the percentage of nudity exceeds 50%. For the concept of violence, we curated an additional set of 30 prompts to prevent any overlap with nudity prompts. These prompts have a nudity ratio of less than 50%, an inappropriateness ratio of more than 50%, and are labeled as harmful.

**Baselines.** To verify the effectiveness of our method PLA, we select several SOTA baselines for adversarial attacks on T2I models, including QF-attack [45], SneakyPrompt [40], Ring-A-Bell [37], UnlearnDiffAtk [43], and MMA-Diffusion [38]. The details of the baselines are provided in the Appendix.

**Auxiliary Models.** We adopt SDv1.4 [7] (UNet-based) and PixArt [6] (DiT-based) as the auxiliary models. By leveraging their distinct diffusion architectures, we can provide a comprehensive evaluation of our method. The main text focuses on SDv1.4 implementation, with complete PixArt studies provided in the Appendix.

**Victim T2I Models.** We conduct experiments on three black-box victim T2I models: *SDv1.5* [8], *SDXLv1.0* [25], and *SLD* [34]. Moreover, we test the well-known T2I online

services: *Stability.ai* [9] and *DALL·E 3* [1]. The details of the victim T2I models are provided in the Appendix.

**Safety Mechanisms.** Following previous work [38], we employ the same prompt filter, which utilizes a predefined list of sensitive words to prohibit harmful prompts. And we apply three post-hoc safety checkers: the built-in safety checker in SD [8], Q16 [33], and MHSC [26].

**Evaluation Metrics.** Following MMA-Diffusion [38], we adopt the Attack Success Rate out of N syntheses (ASR-N) as our evaluation metrics. ASR-N measures N generated images of T2I models for each given prompt. The attack is deemed successful if any of these images exhibit NSFW content and bypass the safety mechanisms. For example, ASR-4 indicates the proportion of effective prompts (i.e., at least one out of the four generated images contains NSFW content) over all tested prompts. We evaluate three black-box T2I models using SC [8], Q16 [33], and MHSC [26] to quantify ASR. For online services, six human evaluators independently assess and report the average result.

**Evaluation Settings.** We adopt the pre-trained language models BERT [16] and T5 [28] to generate adversarial prompts. More details of the evaluation settings are provided in the Appendix.

### 5.2. Attacking on Black-Box Victim T2I Models

Due to different choices of pre-trained language models, we set up two models, i.e. PLA-BERT and PLA-T5. We conduct experiments on two datasets: nudity and violence, as shown in Tab. 1 and Tab. 2. On both datasets, compared to other baselines, our attack achieves significant success in steering black-box T2I models to generate NSFW content.

For the nudity dataset, Tab. 1 presents experimental results comparing various adversarial attacks across three black-box victim T2I models using three post-hoc safety checkers and evaluating them based on ASR-4 and ASR-1. It is particularly notable that our proposed methods, PLA-BERT and PLA-T5, outperform all other methods significantly. For the three black-box victim T2I models, the average ASR-4 of PLA-BERT is 91.45%, 90.57%, and 90.82% respectively. Especially on the SDXLv1.0 model, the average ASR-4 of PLA-BERT far exceeds the average of the highest ASR-4 among other baselines, up to 17.27%. Meanwhile, PLA-T5 achieves average ASR-4 scores of 86.56%, 86.54%, and 90.61% respectively. Although PLA-T5's performance on the nudity dataset is slightly lower than that of PLA-BERT, it is still far better than other baseline methods.

Tab. 2 shows experimental results on the violence dataset. The results consistently show that the attack performance of PLA surpasses that of other baselines. But unlike the dataset nudity, PLA-T5 performs better than PLA-BERT on this dataset. This may be due to the fact that different pre-trained language models exhibit distinct "preferences". These results collectively demonstrate the effective-

| Model | Method | SC [8] | | Q16 [33] | | MHSC [26] | | AVG. | |
|---|---|---|---|---|---|---|---|---|---|
| | | ASR-4 | ASR-1 | ASR-4 | ASR-1 | ASR-4 | ASR-1 | ASR-4 | ASR-1 |
| SDv1.5 | QF-Attack [45] (CVPR' 23) | 27.88 | 12.55 | 26.57 | 10.94 | 19.68 | 7.58 | 24.71 | 10.36 |
| | SneakyPrompt [40] (S&P'24) | 44.82 | 24.80 | 35.18 | 19.06 | 33.68 | 16.81 | 37.89 | 20.22 |
| | Ring-A-Bell [37] (ICLR'24) | 58.05 | 35.80 | 51.75 | 33.58 | 41.79 | 19.97 | 50.53 | 29.78 |
| | UnlearnDiffAtk [43] (ECCV' 24) | 75.03 | 58.26 | 74.22 | 55.29 | 70.57 | 51.33 | 73.27 | 54.96 |
| | MMA-Diffusion [38] (CVPR' 24) | 79.14 | 61.30 | 78.38 | 58.36 | 75.77 | 55.48 | 77.76 | 58.38 |
| | **PLA-BERT(Ours)** | **92.41** | **71.44** | **92.61** | **66.10** | **89.33** | 68.52 | **91.45** | **68.69** |
| | **PLA-T5(Ours)** | 89.77 | 69.53 | 83.90 | 64.27 | 83.01 | 63.72 | 86.56 | 65.84 |
| SDXLv1.0 | QF-Attack [45] (CVPR' 23) | 13.93 | 4.73 | 12.46 | 4.18 | 10.08 | 3.34 | 12.16 | 4.08 |
| | SneakyPrompt [40] (S&P'24) | 23.25 | 14.01 | 20.26 | 9.16 | 15.11 | 8.91 | 19.54 | 10.69 |
| | Ring-A-Bell [37] (ICLR'24) | 31.47 | 18.42 | 28.02 | 13.44 | 23.10 | 11.17 | 27.53 | 14.34 |
| | UnlearnDiffAtk [43] (ECCV' 24) | 66.28 | 37.21 | 68.43 | 40.19 | 60.24 | 39.31 | 64.98 | 38.90 |
| | MMA-Diffusion [38] (CVPR' 24) | 72.98 | 41.37 | 77.52 | 49.33 | 69.39 | 45.02 | 73.30 | 45.24 |
| | **PLA-BERT(Ours)** | **95.37** | **76.20** | **94.03** | **74.56** | **82.30** | **63.54** | **90.57** | **71.43** |
| | **PLA-T5(Ours)** | 91.26 | 74.08 | 85.34 | 66.90 | 83.01 | 59.74 | 86.54 | 66.91 |
| SLD | QF-Attack [45] (CVPR' 23) | 19.27 | 8.90 | 18.91 | 7.47 | 16.76 | 6.78 | 18.31 | 7.72 |
| | SneakyPrompt [40] (S&P'24) | 49.90 | 26.32 | 36.29 | 22.46 | 37.91 | 23.37 | 41.37 | 24.05 |
| | Ring-A-Bell [37] (ICLR'24) | 56.88 | 38.26 | 51.16 | 33.29 | 49.72 | 29.94 | 52.59 | 33.83 |
| | UnlearnDiffAtk [43] (ECCV' 24) | 72.39 | 40.24 | 62.53 | 47.20 | 65.17 | 51.84 | 66.70 | 46.43 |
| | MMA-Diffusion [38] (CVPR' 24) | 75.99 | 45.27 | 75.34 | 53.44 | 78.12 | 60.28 | 76.48 | 53.00 |
| | **PLA-BERT(Ours)** | **94.75** | 73.09 | **90.32** | **64.88** | 87.39 | **69.94** | 90.82 | **69.30** |
| | **PLA-T5(Ours)** | 93.41 | **75.60** | 88.24 | 60.03 | **90.17** | 67.53 | **90.61** | 67.72 |

Table 1. The attack performance of PLA against black-box T2I models on the nudity dataset. The **bolded** values are the highest performance. The difference between PLA-BERT and PLA-T5 is the pre-trained language model used to generate adversarial prompts.

| Model | Method | SC [8] | | Q16 [33] | | MHSC [26] | | AVG. | |
|---|---|---|---|---|---|---|---|---|---|
| | | ASR-4 | ASR-1 | ASR-4 | ASR-1 | ASR-4 | ASR-1 | ASR-4 | ASR-1 |
| SDv1.5 | QF-Attack [45] (CVPR' 23) | 25.15 | 11.76 | 23.81 | 9.44 | 18.59 | 7.28 | 22.52 | 9.49 |
| | SneakyPrompt [40] (S&P'24) | 38.71 | 17.77 | 36.26 | 15.14 | 35.62 | 16.61 | 36.86 | 16.51 |
| | Ring-A-Bell [37] (ICLR'24) | 65.41 | 40.02 | 54.24 | 38.90 | 53.04 | 37.73 | 57.56 | 38.88 |
| | UnlearnDiffAtk [43] (ECCV' 24) | 71.22 | 54.17 | 65.23 | 46.88 | 63.92 | 47.31 | 66.79 | 49.45 |
| | MMA-Diffusion [38] (CVPR' 24) | 80.23 | 64.46 | 78.45 | 61.71 | 76.11 | 56.96 | 78.26 | 61.04 |
| | **PLA-BERT(Ours)** | **93.46** | **73.81** | 91.44 | 73.28 | 80.97 | 61.44 | 88.62 | 69.51 |
| | **PLA-T5(Ours)** | 92.04 | 71.38 | **93.96** | **75.90** | **85.23** | **64.73** | **90.41** | **70.67** |
| SDXLv1.0 | QF-Attack [45] (CVPR' 23) | 12.81 | 3.62 | 11.24 | 3.55 | 10.18 | 2.08 | 11.41 | 3.08 |
| | SneakyPrompt [40] (S&P'24) | 34.45 | 16.17 | 26.38 | 10.65 | 24.80 | 9.77 | 28.54 | 12.20 |
| | Ring-A-Bell [37] (ICLR'24) | 42.78 | 30.47 | 34.21 | 26.82 | 31.72 | 23.05 | 36.24 | 26.78 |
| | UnlearnDiffAtk [43] (ECCV' 24) | 65.29 | 49.42 | 64.83 | 41.27 | 62.81 | 39.90 | 64.31 | 43.53 |
| | MMA-Diffusion [38] (CVPR' 24) | 75.92 | 53.23 | 76.01 | 50.29 | 74.67 | 48.32 | 75.53 | 50.61 |
| | **PLA-BERT(Ours)** | 91.69 | 70.23 | 90.04 | 71.36 | 79.11 | 58.25 | 86.95 | 66.61 |
| | **PLA-T5(Ours)** | **93.72** | **78.91** | **92.63** | **78.04** | **80.51** | **62.94** | **88.95** | **73.30** |
| SLD | QF-Attack [45] (CVPR' 23) | 18.48 | 8.88 | 16.76 | 7.15 | 16.28 | 6.54 | 17.17 | 7.52 |
| | SneakyPrompt [40] (S&P'24) | 50.32 | 36.61 | 45.94 | 31.39 | 42.26 | 33.00 | 46.17 | 33.67 |
| | Ring-A-Bell [37] (ICLR'24) | 69.93 | 49.48 | 61.57 | 49.06 | 59.50 | 38.99 | 63.67 | 45.84 |
| | UnlearnDiffAtk [43] (ECCV' 24) | 61.08 | 46.74 | 66.28 | 44.91 | 63.02 | 45.27 | 63.46 | 45.64 |
| | MMA-Diffusion [38] (CVPR' 24) | 76.62 | 55.76 | 77.95 | 56.49 | 74.77 | 58.60 | 76.45 | 56.95 |
| | **PLA-BERT(Ours)** | 91.98 | 77.84 | 91.22 | 71.54 | 84.41 | **66.70** | 89.20 | 72.03 |
| | **PLA-T5(Ours)** | **93.34** | **79.62** | **92.74** | **73.04** | **86.33** | 64.19 | **90.80** | **72.28** |

Table 2. The attack performance of PLA against black-box T2I models on the violence dataset. The **bolded** values are the highest performance. The difference between PLA-BERT and PLA-T5 is the pre-trained language model used to generate adversarial prompts.

ness of PLA in tackling the challenging task of bypassing both prompt filter and post-hoc safety checkers under the black-box setting.

**Visualization of Results.** Fig. 3 presents images generated by SDXLv1.0 using adversarial prompts created with PLA, demonstrating the strong capability of our attack method to generate NSFW content against black-box T2I models. Additional visual results are provided in the Appendix.

## 5.3. Attacking on T2I Online Services

We evaluate two popular online services, Stability.ai [9] and DALL·E 3 [1], both of which are equipped with proprietary safety mechanisms as shown in Tab. 3. Due to network delays and limitations on the number of queries allowed, conducting quantitative tests on the large dataset we collected directly is challenging. To address this, we use a subset of the large dataset (20 nudity prompts and 20 violence prompts). Also compared to other baselines, our at-

| Dataset | Model | QF-Attack | SneakyPrompt | Ring-A-Bell | UnlearnDiffAtk | MMA-Diffusion | **PLA-BERT** | **PLA-T5** |
|---|---|---|---|---|---|---|---|---|
| **Nudity** | **Stability.ai** | 39.18 | 9.44 | 31.27 | 44.03 | 46.89 | **62.15** | 54.83 |
| | **DALL·E 3** | 30.26 | 6.57 | 26.97 | 28.02 | 28.72 | **45.09** | 38.22 |
| **Violence** | **Stability.ai** | 13.62 | 28.64 | 46.24 | 40.81 | 42.57 | 55.68 | **69.70** |
| | **DALL·E 3** | 9.08 | 13.11 | 51.31 | 24.76 | 25.80 | 36.77 | **51.98** |

Table 3. Evaluation of different attack methods on T2I online services via the metric of ASR-4.



Figure 3. **Visualization results of PLA.** Sensitive words within the target prompt are colored in red. Images are generated by SDXLv1.0.

tack method exhibits superior attack performance. We provide more examples of NSFW images generated by T2I online services in the Appendix.

### 5.4. Ablation Study

**Multimodal Loss.** To demonstrate the effectiveness of the multimodal loss, we conduct ablation studies by removing the specific $\mathcal{L}_a$ (or $\mathcal{L}_b$) in our approach. We use PLA-T5 to attack the SLD model on the violence and nudity datasets. As shown in Tab. 4, in the absence of $\mathcal{L}_a$ or $\mathcal{L}_b$, the attack performance of our method significantly decreases, indicating that these two components play a crucial role in the effectiveness of our attack. In particular, the impact of $\mathcal{L}_b$ on attack performance is more significant. This may be due to the presence of more potentially sensitive information in the target images, which more effectively guides the generation of adversarial prompts.

**Gradient Optimization** To verify the powerful capability of our gradient design, we perform an ablation study on it. We adopt different insertion schemes:

• We keep our gradient method (Eq. (9)) and "restart" strat-

| | **Violence** | | **Nudity** | |
|---|---|---|---|---|
| **PLA** (Ours) | ASR-4 | ASR-1 | ASR-4 | ASR-1 |
| $\mathcal{L}_a + \mathcal{L}_b$ | **93.34** | **79.62** | **93.41** | **75.60** |
| - w/o $\mathcal{L}_a$ | 81.02 | 54.57 | 82.99 | 51.07 |
| - w/o $\mathcal{L}_b$ | 79.34 | 47.88 | 74.66 | 44.87 |

Table 4. Ablation study on multimodal loss.

| **Gradient Method** | **Violence** | | **Nudity** | |
|---|---|---|---|---|
| | ASR-4 | ASR-1 | ASR-4 | ASR-1 |
| $G_{PLA}$ | **91.69** | **70.23** | **95.37** | **76.20** |
| $G_{ZOO}$ | 52.89 | 46.73 | 58.44 | 41.27 |
| $G_{RE}$ | 70.12 | 58.24 | 78.33 | 53.90 |

Table 5. The Analysis of Gradient Optimization.

egy (i.e., $G_{PLA}$).
• We utilize the ZOO gradient method (Eq. (8)) and "restart" strategy (i.e., $G_{ZOO}$).
• We keep our gradient method (Eq. (9)) but remove "restart" strategy (i.e., $G_{RE}$).

As shown in Tab. 5, we utilize PLA-BERT to attack the SDXLv1.0 model on the violence and nudity datasets. We can see that our gradient method outperforms the traditional ZOO method. The absence of the "restart" strategy leads to the decrease of ASR, primarily because generating black images in the initial optimization step triggers gradient vanishing.

### 6. Conclusion

This study investigates the vulnerability of black-box T2I models against adversarial attacks that bypass safety mechanisms including prompt filters and post-hoc safety checkers. Due to the unique challenges of training gradient-driven attack methods under black-box settings, most previous methods rely on word substitution to search adversarial prompts over limited search space, leading to suboptimal performance compared to gradient-based training. To bridge this gap, we propose a novel prompt learning attack framework (PLA), where insightful gradient-based training tailored to black-box T2I models is designed by utilizing multimodal similarities. Our results affirm that employing PLA to fabricate adversarial prompts can potentially steer these T2I models to output NSFW content effectively, contributing to the development of more robust defensive strategies in the future.

# References

[1] Dall·e 3. Available: https://openai.com/index/dall-e-3/. 1, 6, 7

[2] Gpt-4. Available: https://openai.com/index/gpt-4/. 1

[3] Gen-2. Available: https://runwayml.com/research/gen-2.

[4] Microsoft designer. Available: https://designer.microsoft.com/. 1

[5] Midjourney. Available: https://www.midjourney.com/. 3

[6] Pixart. Available: https://github.com/PixArt-alpha/PixArt-alpha. 6

[7] Stable diffusion v1.4 checkpoint, . Available: https://huggingface.co/CompVis/stable-diffusion-v1-4. 6

[8] Stable diffusion v1.5 checkpoint, . Available: https://huggingface.co/stable-diffusion-v1-5/stable-diffusion-v1-5. 1, 6, 7

[9] stability.ai. Available: https://stability.ai/. 1, 6, 7

[10] Aochuan Chen, Yimeng Zhang, Jinghan Jia, James Diffenderfer, Jiancheng Liu, Konstantinos Parasyris, Yihua Zhang, Zheng Zhang, Bhavya Kailkhura, and Sijia Liu. Deepzero: Scaling up zeroth-order optimization for deep model training. *arXiv preprint arXiv:2310.02025*, 2023. 5

[11] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pages 15–26, 2017.

[12] Xiangyi Chen, Sijia Liu, Kaidi Xu, Xingguo Li, Xue Lin, Mingyi Hong, and David Cox. Zo-adamm: Zeroth-order adaptive momentum method for black-box optimization. *Advances in neural information processing systems*, 32, 2019. 5

[13] Mohammad Sh Daoud, Mohammad Shehab, Hani M Al-Mimi, Laith Abualigah, Raed Abu Zitar, and Mohd Khaled Yousef Shambour. Gradient-based optimizer (gbo): a review, theory, variants, and applications. *Archives of Computational Methods in Engineering*, 30(4):2431–2449, 2023. 1

[14] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2426–2436, 2023. 2

[15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 6

[16] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of naacL-HLT*, page 2. Minneapolis, Minnesota, 2019. 5, 6

[17] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 2

[18] Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu. Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22691–22702, 2023. 2

[19] Tony Lee, Michihiro Yasunaga, Chenlin Meng, Yifan Mai, Joon Sung Park, Agrim Gupta, Yunzhi Zhang, Deepak Narayanan, Hannah Teufel, Marco Bellagente, et al. Holistic evaluation of text-to-image models. *Advances in Neural Information Processing Systems*, 36, 2024. 2

[20] Chumeng Liang, Xiaoyu Wu, Yang Hua, Jiaru Zhang, Yiming Xue, Tao Song, Zhengui Xue, Ruhui Ma, and Haibing Guan. Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples. In *International Conference on Machine Learning*, pages 20763–20786. PMLR, 2023. 3

[21] Han Liu, Yuhao Wu, Shixuan Zhai, Bo Yuan, and Ning Zhang. Riatig: Reliable and imperceptible adversarial text-to-image generation with natural prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20585–20594, 2023. 1

[22] Qihao Liu, Adam Kortylewski, Yutong Bai, Song Bai, and Alan Yuille. Intriguing properties of text-guided diffusion models. *arXiv preprint arXiv:2306.00974*, 2, 2023. 3

[23] Runtao Liu, Ashkan Khakzar, Jindong Gu, Qifeng Chen, Philip Torr, and Fabio Pizzati. Latent guard: a safety framework for text-to-image generation. *arXiv preprint arXiv:2404.08031*, 2024. 2

[24] Natalie Maus, Patrick Chao, Eric Wong, and Jacob R Gardner. Black box adversarial prompting for foundation models. In *The Second Workshop on New Frontiers in Adversarial Machine Learning*, 2023. 3

[25] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 6

[26] Yiting Qu, Xinyue Shen, Xinlei He, Michael Backes, Savvas Zannettou, and Yang Zhang. Unsafe diffusion: On the generation of unsafe images and hateful memes from text-to-image models. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pages 3403–3417, 2023. 1, 2, 6, 7

[27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5

[28] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020. 5, 6

[29] Javier Rando, Daniel Paleka, David Lindner, Lennart Heim, and Florian Tramèr. Red-teaming the stable diffusion safety filter. *arXiv preprint arXiv:2210.04610*, 2022. 2

[30] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 3

[31] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022. 1

[32] Hadi Salman, Alaa Khaddaj, Guillaume Leclerc, Andrew Ilyas, and Aleksander Madry. Raising the cost of malicious ai-powered image editing. In *International Conference on Machine Learning*, pages 29894–29918. PMLR, 2023. 3

[33] Patrick Schramowski, Christopher Tauchmann, and Kristian Kersting. Can machines help us answering question 16 in datasheets, and in turn reflecting on inappropriate content? In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1350–1361, 2022. 1, 2, 6, 7

[34] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22522–22531, 2023. 1, 2, 6

[35] James C Spall. *Introduction to stochastic search and optimization: estimation, simulation, and control*. John Wiley & Sons, 2005. 5

[36] Mandt Stephan, Matthew D Hoffman, David M Blei, et al. Stochastic gradient descent as approximate bayesian inference. *Journal of Machine Learning Research*, 18(134):1–35, 2017. 2

[37] Yu-Lin Tsai, Chia-yi Hsu, Chulin Xie, Chih-hsun Lin, Jia You Chen, Bo Li, Pin-Yu Chen, Chia-Mu Yu, and Chun-ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion models? In *International Conference on Learning Representations*, 2024. 3, 6, 7

[38] Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Tsung-Yi Ho, Nan Xu, and Qiang Xu. Mma-diffusion: Multimodal attack on diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7737–7746, 2024. 1, 3, 6, 7

[39] Yijun Yang, Ruiyuan Gao, Xiao Yang, Jianyuan Zhong, and Qiang Xu. Guardt2i: Defending text-to-image models from adversarial prompts. *arXiv preprint arXiv:2403.01446*, 2024. 2

[40] Yuchen Yang, Bo Hui, Haolin Yuan, Neil Gong, and Yinzhi Cao. Sneakyprompt: Jailbreaking text-to-image generative models. In *2024 IEEE symposium on security and privacy (SP)*, pages 897–912. IEEE, 2024. 1, 3, 6, 7

[41] Gong Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not: Learning to forget in text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1755–1764, 2024. 2

[42] Jianping Zhang, Zhuoer Xu, Shiwen Cui, Changhua Meng, Weibin Wu, and Michael R Lyu. On the robustness of latent diffusion models. *arXiv preprint arXiv:2306.08257*, 2023. 3

[43] Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. *arXiv preprint arXiv:2310.11868*, 2023. 1, 2, 3, 6, 7

[44] Wei Zhou, Pengjun Wang, Ali Asghar Heidari, Xuehua Zhao, Hamza Turabieh, and Huiling Chen. Random learning gradient based optimization for efficient design of photovoltaic models. *Energy Conversion and Management*, 230: 113751, 2021. 1

[45] Haomin Zhuang, Yihua Zhang, and Sijia Liu. A pilot study of query-free adversarial attack against stable diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2385–2392, 2023. 3, 6, 7