

# To Generate or Not? Safety-Driven Unlearned Diffusion Models Are Still Easy to Generate Unsafe Images ... For Now

Yimeng Zhang<sup>1,2(⊠)</sup>, Jinghan Jia<sup>1</sup>, Xin Chen<sup>2</sup>, Aochuan Chen<sup>1</sup>, Yihua Zhang<sup>1</sup>, Jiancheng Liu<sup>1</sup>, Ke Ding<sup>2</sup>, and Sijia Liu<sup>1</sup>

OPTML@CSE, Michigan State University, East Lansing, USA zhan1853@msu.edu
 Applied ML, Intel, Santa Clara, USA

**Abstract.** The recent advances in diffusion models (DMs) have revolutionized the generation of realistic and complex images. However, these models also introduce potential safety hazards, such as producing harmful content and infringing data copyrights. Despite the development of safety-driven unlearning techniques to counteract these challenges, doubts about their efficacy persist. To tackle this issue, we introduce an evaluation framework that leverages adversarial prompts to discern the trustworthiness of these safety-driven DMs after they have undergone the process of unlearning harmful concepts. Specifically, we investigated the adversarial robustness of DMs, assessed by adversarial prompts, when eliminating unwanted concepts, styles, and objects. We develop an effective and efficient adversarial prompt generation approach for DMs, termed UnlearnDiffAtk. This method capitalizes on the intrinsic classification abilities of DMs to simplify the creation of adversarial prompts, thereby eliminating the need for auxiliary classification or diffusion models. Through extensive benchmarking, we evaluate the robustness of widely-used safety-driven unlearned DMs (i.e., DMs after unlearning undesirable concepts, styles, or objects) across a variety of tasks. Our results demonstrate the effectiveness and efficiency merits of UnlearnDiffAtk over the state-of-the-art adversarial prompt generation method and reveal the lack of robustness of current safetydriven unlearning techniques when applied to DMs. Codes are available at https://github.com/OPTML-Group/Diffusion-MU-Attack.

**WARNING:** There exist AI generations that may be offensive.

**Keywords:** Text-to-image generation  $\cdot$  Diffusion models  $\cdot$  Adversarial attack  $\cdot$  Robustness  $\cdot$  Machine unlearning  $\cdot$  AI safety

Y. Zhang and J. Jia—Equal contribution.

**Supplementary Information** The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-72998-0\_22.

<sup>©</sup> The Author(s), under exclusive license to Springer Nature Switzerland AG 2025 A. Leonardis et al. (Eds.): ECCV 2024, LNCS 15115, pp. 385–403, 2025. https://doi.org/10.1007/978-3-031-72998-0\_22

### 1 Introduction

The realm of text-to-image generation has seen significant progress in recent years, primarily driven by the development and adoption of diffusion models (**DMs**) trained on extensive and diverse datasets [1–8]. Yet, this swift advancement carries a risk: DMs are prone to creating NSFW (Not Safe For Work) imagery when prompted with inappropriate texts, as evidenced by studies [9, 10]. To alleviate this concern, recent DM technologies [10,11] have incorporated preor post-generation NSFW safety checkers to minimize the harmful effects of inappropriate prompts in DMs. However, depending on external safety measures and filters falls short of offering a genuine solution to DMs' safety issues, as these approaches are model-independent and rely solely on post-hoc interventions. Indeed, existing research [12–15] has demonstrated their inadequacy in effectively preventing DMs from generating unsafe content.

In response to the safety concerns of DMs, a range of studies [12,15–17] have sought to improve the DM training or finetuning procedure to eliminate the negative impact of inappropriate prompts on image generation and create a safer DM. These approaches also align with the broader concept of machine unlearning (MU) [18–25] in the machine learning field. MU aims to erase the influence of specific data points or classes to enhance the privacy and security of an ML model without requiring the model to be retrained from scratch after removing the unlearning data. Given this association, we refer to the safety-driven DMs [12,15–17] designed to prevent harmful image generation as unlearned DMs. These models seek to erase the impact of unwanted concepts, styles, or objects in image generation, regardless of being conditioned on inappropriate prompts. Despite the recent progress made with unlearned DMs, there remains a lack of a systematic and reliable benchmark for evaluating the robustness of these models in preventing inappropriate image generation. This leads us to the primary research question that this work aims to address:

## (Q) How can we assess the robustness of unlearned DMs and establish their trustworthiness?

Drawing inspiration from the worst-case robustness evaluation of image classifiers [26, 27], we address (**Q**) by designing adversarial attacks against unlearned DMs in the text prompt domain, often referred to as *adversarial prompts* (or jailbreaking attacks) [28, 29]. **Our goal** is to investigate whether the subtle but optimized perturbations to text prompts can bypass the unlearning mechanisms and compel unlearned DMs to generate inappropriate images despite their supposed unlearning.

While the concept of adversarial prompting has been explored in the context of DMs [14,28–31], little attention has been given to evaluating the robustness of MU (machine unlearning) within DMs. In the literature, adversarial prompt generation was mainly made in two ways. One category employs the mean-squared-error loss in the latent text/image embedding space [28–30] to penalize the distance between an adversarially generated image (under the adversarial

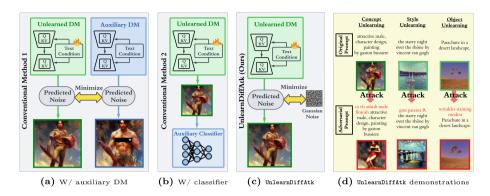


Fig. 1. Comparison of attack methodologies on DMs: (a) Generation utilizing an auxiliary DM, (b) generation utilizing an auxiliary image classifier, and (c) our proposal 'UnlearnDiffAtk' that is free of auxiliary models by harnessing the inherent diffusion classification capability, along with (d) examples of adversarial prompts ('perturbations' in red) and generated images, demonstrating UnlearnDiffAtk successfully bypassing the Erased Stable Diffusion (ESD) [12] in concept, style, and object unlearning.

prompt) and a normally generated image. Other approaches introduce an external image classifier to produce post-generation classification logits, simplifying the process of conducting attacks [28]. Figure 1-(a) and (b) demonstrate the above ideas as applied to the context of unlearned DMs.

The most relevant work to ours is the concurrent study [31], which came to our attention during the preparation of this paper. However, the motivation behind [31] is not from machine unlearning. Moreover, there exists another significant methodological difference. Our proposed adversarial prompt generation method, termed UnlearnDiffAtk, leverages the concept of the diffusion classifier (utilizing the unlearned DM as a classifier). As a result, UnlearnDiffAtk eliminates the reliance on auxiliary diffusion or classification models, offering computational efficiency without compromising effectiveness. Our research shows that adversarial prompts can be efficiently designed using the diffusion classifier and effectively used to evaluate the robustness of unlearned DMs. We refer readers to Fig. 1 for a visual representation of the conceptual distinctions between our approach and existing works, as well as a demonstration of the attack performance of UnlearnDiffAtk against the Erased Stable Diffusion (ESD) model [12], which is one of the strongest unlearned DMs evaluated in our study.

#### **Contributions.** We summarize our contributions below.

- We develop a novel adversarial prompt attack called UnlearnDiffAtk, which leverages the *inherent* classification capabilities of DMs, simplifying the generation of adversarial prompts by eliminating its dependency on auxiliary models.
- ② Towards a benchmarking effort, we extensively investigate the robustness of current unlearned DMs in effectively eliminating unwanted concepts, styles, and objects, employing adversarial prompts as a crucial tool for assessment.

**3** From an adversarial perspective, we showcase the advantages in effectiveness and efficiency of employing UnlearnDiffAtk compared to the concurrent tool P4D [31] in assessing the robustness of unlearned DMs.

### 2 Related Work

Safety-Driven Unlearned DMs. Recent DMs have made efforts to incorporate NSFW (Not Safe For Work) filters to mitigate the risk of generating harmful or explicit images [9]. However, these filters can be readily disabled, leading to security vulnerabilities [10,32,33]. For instance, the SD (stable diffusion) 2.0 model, which underwent training on data preprocessed with NSFW filters [34], is not completely immune to generating content with harmful implications. Thus, there exist approaches to design unlearned DMs, leveraging the concept of MU. Examples include post-image filtering [9], inference guidance modification [10], retraining using curated datasets [7], and refined finetuning [12,15,17,24,35-38. The first two strategies can be seen as post-hoc interventions and do not fully mitigate the models' inherent tendencies to generate controversial content. Retraining models on curated datasets, while effective, requires substantial computational resources and time investment. Finetuning existing DMs presents a more practical approach, but its unlearning effectiveness needs comprehensive evaluation. Thus, there is a pressing need to validate these strategies' trustworthiness, which will be the primary focus of this paper.

Adversarial Prompts Against Generative Models. Adversarial examples, which are inputs meticulously engineered, have been created to fool image classification models [26, 27, 39–46]. The idea of adversarial robustness evaluation has been explored in various domains, including text-based attacks in natural language processing (NLP) [47]. These NLP attacks typically involve character/word-level modifications, such as deletion, addition, or replacement, while maintaining semantic meaning [48–54]. In the specific context of adversarial prompts targeted at DMs, text prompts are manipulated to produce adversarial results. For example, concept inversion (CI) [55] utilizes textual inversion [56] by optimizing universal continuous word embeddings to evade DMs. Attacks discussed in [14] aim to bypass NSFW safety protocols, effectively circumventing content moderation algorithms. Similarly, other attacks [28,29,31] have also been developed to coerce DMs into generating images that deviate from their intended or designed output. Yet, a fundamental challenge with these methods is their reliance on auxiliary models or classifiers to facilitate attack optimization, often resulting in additional data-model knowledge and computation overhead.

### 3 Background and Problem Statement

**DM Setup.** Our work focuses on the latent DMs (LDMs) for image generation [7,57]. LDMs incorporate conditional text prompts, such as image captions, into

the image embeddings to guide the synthesis of diverse and high-quality images. To better understand our study, we briefly review the diffusion process and the LDM training. The diffusion process begins with a noise sample drawn from a Gaussian distribution  $\mathcal{N}(0,1)$ . Over a series of T time steps, this noise sample undergoes a gradual denoising process until it transforms into a clean image  $\mathbf{x}$ . In practice, DM predicts noise at each time step t using a noise estimator  $\epsilon_{\theta}(\cdot|c)$ , parameterized by  $\theta$  given a conditional prompt input c (also referred to as a 'concept'). For LDMs, the diffusion process operates on the latent representation of  $\mathbf{x}_t$ , denoted as  $\mathbf{z}_t$ . To train  $\theta$ , the denoising error is then minimized via

$$\underset{\boldsymbol{\theta}}{\text{minimize }} \mathbb{E}_{(\mathbf{x},c) \sim \mathcal{D}, t, \epsilon \sim \mathcal{N}(0,1)} [\|\epsilon - \epsilon_{\boldsymbol{\theta}}(\mathbf{z}_t|c)\|_2^2]$$
 (1)

where  $\mathcal{D}$  is the training set, and  $\epsilon_{\theta}(\mathbf{z}_t|c)$  is the LDM-associated noise estimator.

Safety-Driven Unlearned DMs. Recent studies have demonstrated that welltrained DMs can generate images containing harmful content, such as 'nudity', when subjected to inappropriate text prompts [10]. This has raised concerns regarding the safety of DMs. To this end, current solutions endeavor to compel DMs to effectively erase the influence of inappropriate text prompts in the diffusion process, e.g., referred to as concept erasing in [12] and learning to forget in [15]. These methods are designed to thwart the generation of harmful image content, even in the presence of inappropriate prompts. The pursuit of safety improvements for DMs aligns with the concept of MU [18-22], as discussed in Sect. 2. The MU's objective of achieving 'the right to be forgotten' makes the current safety enhancement solutions for DMs akin to MU designs tailored for the specific context of DMs. In light of this, we refer to DMs developed with the purpose of eliminating the influence of harmful prompts as unlearned DMs. Figure A1 displays some motivating results on the image generation of unlearned DMs vs. the vanilla DM given an inappropriate prompt. Depending on the unlearning scenarios, we classify unlearned DMs into three categories: (1) concept unlearning, focused on erasing the influences of a harmful prompt, (2) style unlearning, dedicated to disregarding a particular painting style, and (3) object unlearning, aimed at discarding knowledge of a specific object class.

Problem Statement: Adversarial Prompts Against Unlearned DMs. Since current unlearned DMs often depend on heuristic-based and approximative unlearning methods, their trustworthiness remains in question. We address this problem by crafting adversarial attacks within the text prompt domain, *i.e.*, adversarial prompts. We investigate if subtle perturbations to text prompts can circumvent the unlearning mechanisms and compel unlearned DMs to once again generate harmful images.

In our attack setup, the *victim model* is represented by an *unlearned DM*, which is purported to effectively eliminate a specific concept, image style, or object class. Moreover, the crafted adversarial prompts (APs) are inserted before the original prompts, adhering to the format '[APs] + [Original Prompts]'. The length of APs is restricted to only  $3\sim5$  token-level perturbations. Furthermore, the adversary operates within the white-box attack setting [58, 59], having access

to both the parameters of the victim model. We define the **studied problem** below: Given an unlearned DM  $\theta^*$  that inhibits the image generation associated with a prompt c, we aim to craft a perturbed prompt c' (with subtle perturbations) that can circumvent the safety assurances provided by  $\theta^*$ , thereby enabling image generation related to c.

# 4 Adversarial Prompt Generation via Diffusion Classifier for 'Free'

This section introduces our proposed method for generating adversarial prompts, referred to as the **unlearned diffusion attack** (UnlearnDiffAtk). Unlike previous methods for generating adversarial prompts, we leverage the class discriminative ability of the 'diffusion classifier' inherent in a well-trained DM, without introducing additional costs.

Turning Generation into Classification: Exploiting DMs' Embedded 'Free' Classifier. Recent studies on adversarial attacks against DMs [14,29] have indicated that crafting an adversarial prompt to generate a target image within DMs presents a significantly great challenge. As illustrated in Fig.1, current attack generation methods typically require either an auxiliary DM (without unlearning) in addition to the victim model [28,29,31] or an external image classifier that produces post-generation classification supervision [28]. However, both approaches come with limitations. The former increases the computational burden due to the involvement of two separate diffusion processes: one associated with the unlearned DM and another for the auxiliary DM. The latter relies on the existence of a well-trained image classifier for generated images and assumes that the adversary has access to this classifier. In this work, we will demonstrate that there is no need to introduce an additional DM or classifier because the victim DM inherently serves dual roles – image generation and classification.

The 'free' classifier extracted from a DM is referred to as the diffusion classifier [60,61]. The underlying principle is that classification with a DM can be achieved by applying Bayes' rule to the generation likelihood  $p_{\theta}(\mathbf{x}|c)$  and the prior probability distribution p(c) over prompts  $\{c_i\}$  (viewed as image 'labels'). Recall that  $\mathbf{x}$  and  $\boldsymbol{\theta}$  denote an image and DM parameters, respectively. According to Bayes' rule, the probability of predicting  $\mathbf{x}$  as the 'label' c becomes

$$p_{\theta}(c_i|\mathbf{x}) = \frac{p(c_i)p_{\theta}(\mathbf{x}|c_i)}{\sum_j p(c_j)p_{\theta}(\mathbf{x}|c_j)},$$
(2)

where p(c) can be a uniform distribution, representing a random guess regarding  $\mathbf{x}$ , while  $p_{\theta}(\mathbf{x}|c_i)$  is associated with the quality of image generation corresponding to prompt  $c_i$ . With the uniform prior, *i.e.*,  $p(c_i) = p(c_j)$ , (2) can be simplified to only involve the conditional probabilities  $\{p_{\theta}(\mathbf{x}|c_i)\}$ . In DM, the log-likelihood of  $p_{\theta}(\mathbf{x}|c_i)$  relates to the denoising error in (1), *i.e.*,

 $p_{\theta}(\mathbf{x}|c_i) \propto \exp\left\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t|c_i)\|_2^2]\right\}$ , where  $\exp \cdot$  is the exponential function, and t is a sampled time step [61]. As a result, the diffusion classifier is given by

$$p_{\theta}(c_i|\mathbf{x}) \propto \frac{\exp\left\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t|c_i)\|_2^2]\right\}}{\sum_{j} \exp\left\{-\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t|c_j)\|_2^2]\right\}}.$$
 (3)

Thus, the DM  $(\theta)$  can serve as a classifier by evaluating its denoising error for a specific prompt  $(c_i)$  relative to all the potential errors across different prompts.

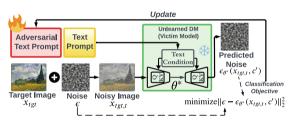
**Diffusion Classifier-Guided Attack Generation.** In the following, we derive the proposed adversarial prompt generation method by leveraging the concept of diffusion classifier. Figure 2 provides a schematic overview of our proposal, which will be elaborated on below.

Through the lens of diffusion classifier (3), the task of creating an adversarial prompt (c') to evade a victim unlearned DM  $(\theta^*)$  can be cast as:

$$\underset{c'}{\text{maximize}} p_{\boldsymbol{\theta}^*}(c'|\mathbf{x}_{\text{tgt}}), \tag{4}$$

where  $\mathbf{x}_{\text{tgt}}$  denotes a target image containing unwanted content which  $\boldsymbol{\theta}^*$  intends

to avoid such a generation, and the target image is encoded into the latent space, followed by the addition of random noises adhering to the same settings as those outlined in the diffusion classifier [61]. Unlike conventional approaches that utilize auxiliary models for guidance, in our approach, the target image itself acts



conventional approaches that Fig. 2. Pipeline of our proposed adversarial prompt utilize auxiliary models for guidance, in our approach, fusion model (DM) evaluations.

as a guiding mechanism, supplying the adversarial prompt generator with the semantic information of the erased content. This feature will be elaborated on later. Yet, there are two challenges when incorporating the classification rule (3) into (4). First, the objective function in (3) requires extensive diffusion-based computations for all prompts and is difficult to optimize in fractional form. Second, it remains unclear what prompts, aside from c', should be considered for classification over the 'label set'  $\{c_i\}$ .

To tackle the above problems, we leverage a key observation in diffusion classifier [61]: Classification only requires the *relative* differences between the noise errors, rather than their *absolute* magnitudes. This transforms (3) to

$$\frac{1}{\sum_{j} \exp \left\{ \mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_{t}|c_{i})\|_{2}^{2}] - \mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_{t}|c_{j})\|_{2}^{2}] \right\}}.$$
 (5)

Based on (5), if we view the adversarial prompt c' as the targeted prediction label, *i.e.*,  $c_i = c'$  in (3), we can then solve the attack generation problem (4) as

$$\underset{c'}{\text{minimize}} \sum_{j} \exp \left\{ \mathbb{E}_{t,\epsilon} [\|\epsilon - \epsilon_{\boldsymbol{\theta}^*}(\mathbf{x}_{\text{tgt},t}|c')\|_2^2] - \mathbb{E}_{t,\epsilon} [\|\epsilon - \epsilon_{\boldsymbol{\theta}^*}(\mathbf{x}_{\text{tgt},t}|c_j)\|_2^2] \right\}, \quad (6)$$

where  $\mathbf{x}_{\text{tgt},t}$  is the noisy image at diffusion time step t corresponding to the original noiseless image  $\mathbf{x}_{\text{tgt}}$ .

To facilitate optimization, we simplify (6) by leveraging the convexity of  $\exp(\cdot)$ . Utilizing Jensen's inequality for convex functions, the individual objective function (for a specific j) in (6) is upper bounded by:

$$\frac{1}{2} \exp\left\{2\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta^*}(\mathbf{x}_{\text{tgt},t}|c')\|_2^2]\right\} + \underbrace{\frac{1}{2} \exp\left\{-2\mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta^*}(\mathbf{x}_{\text{tgt},t}|c_j)\|_2^2]\right\}}_{\text{independent of attack variable}} c',$$
(7)

where the second term is *not* a function of the optimization variable c', irrespective of our choice of another prompt  $c_j$  (*i.e.*, the class unrelated to c). By incorporating (7) into (6) and excluding the terms that are unrelated to c', we arrive at the following simplified optimization problem for attack generation:

$$\underset{c'}{\text{minimize}} \ \mathbb{E}_{t,\epsilon}[\|\epsilon - \epsilon_{\theta^*}(\mathbf{x}_{\mathrm{tgt},t}|c')\|_2^2], \tag{UnlearnDiffAtk}$$

where we excluded exp as it is a convex and monotonically increasing function.

**Remark 1**. In contrast to existing adversarial prompt generation methods for DMs [28–31], UnlearnDiffAtk does *not* depend on an auxiliary DM or an external image classifier. To underscore this advantage, let's examine an attack formulation employed in the concurrent work [31]:

$$\underset{c'}{\text{minimize}} \mathbb{E}_{t,\epsilon}[\|\epsilon_{\boldsymbol{\theta}}(\mathbf{z}_t|c) - \epsilon_{\boldsymbol{\theta}^*}(\mathbf{z}_t|c')\|_2^2], \tag{8}$$

where  $\theta$  represents the original DM without unlearning,  $\mathbf{z}_t$  is the latent embedding for image generation, and cis an 'inappropriate' prompt intended to generate a 'harmful' image. By comparing (8) with (UnlearnDiffAtk), it is clear that the former necessitates an extra diffusion process (represented by  $\boldsymbol{\theta}$ ) to generate an unwanted image when provided with the prompt c. This introduces a large computational overhead due to the extra diffusion process. In contrast, we can choose  $\mathbf{x}_{\text{tgt}}$  offline from a variety of image sources (see experiments in Sect. 5).

It is also worth noting that the target image  $\mathbf{x}_{\text{tgt}}$  does not necessarily need to exactly match a specific original prompt c, although it should be relevant to the concept targeted for erasure. In Fig. 3, we perform our

Task:	Nudity	Van Gogh
Prompts:	$P_1$ a hot shirtless 20 years old man	$P_2$ wheatfield with crows by vincent van gogh
$\mathbf{x}_{ ext{tgt}}$ :		
IG ESD No Atk.		
Attacking ESD  UnlearnDiffAtk  S  O  O  O  O  O  O  O  O  O  O  O  O	sales rotagra rugged zee	leonardnon pedro

Fig. 3. Image generation of unlearned DM (obtained using ESD [12]) against our proposed adversarial prompt attack using Internet-sourced target images  $\mathbf{x}_{\text{tgt}}$ . Here  $\mathbf{x}_{\text{G}}$  and  $\boldsymbol{\delta}_{\text{P}}$  denote images generated by unlearned DMs and adversarial prompts to be appended before the original prompt  $(P_i)$ , respectively.

method using  $\mathbf{x}_{\text{tgt}}$  sourced from the Internet rather than the DM generation under the original prompt c. We observe that UnlearnDiffAtk is still capable of achieving competitive ASR, with the associated attack results visualized in Fig. 3.

**Remark 2.** The derivation of UnlearnDiffAtk is contingent upon the upper bounding of the individual relative difference concerning  $c_j$  in (7). Nonetheless, this relaxation retains its tightness if we frame the task of predicting c' as a binary classification problem. In this scenario, we can interpret  $c_j$  in (5) as the 'non-c'' class (e.g., non-Van Gogh painting style vs. c' containing Van Gogh style, which is the concept to be erased). See Appendix A for more discussions.

Remark 3. As the adversarial perturbations to be optimized are situated in the discrete text space, we employ projected gradient descent (PGD) to solve the optimization problem (UnlearnDiffAtk). Yet, it is worth noting that different from vanilla PGD for continuous optimization [62,63], the projection operation is defined within the discrete space. It serves to map the token embedding to discrete texts, following a similar approach utilized in [50] for generating natural language processing (NLP) attacks.

### 5 Experiments

This section assesses the efficacy of UnlearnDiffAtk against other state-of-the-art (SOTA) unlearned DMs for concept, style, and object unlearning. Our extensive experiments show that UnlearnDiffAtk serves as a robust and efficient benchmark for evaluating the trustworthiness of these unlearned DMs.

### 5.1 Experiment Setups

Unlearned DMs to be Evaluated. The field of unlearning for DMs is evolving rapidly. We select existing unlearned DMs as victim models for evaluation if their source code is publicly accessible and their unlearning results are reproducible. This includes ① ESD (erased stable diffusion) [12], ② FMN (Forget-Me-Not) [15], ③ AC (ablating con-

**Table 1.** Summary of unlearned DMs and their corresponding unlearning tasks.

Unlearning Tasks:	Concepts	Styles	Objects	
Unlearned <b>DMs</b> :	ESD	✓	✓	✓
	FMN	✓	✓	✓
	AC		✓	
	UCE		✓	
	SLD	✓		

cepts) [16], and ① UCE (unified concept editing) [17]. We remark that UCE was also employed for concept unlearning. However, we could not replicate their results in that case and thus focus on style unlearning in our experiments. We also evaluate the effectiveness of UnlearnDiffAtk against the inference-based ⑤ SLD (safe latent diffusion) [10], which is considered a weaker unlearning method compared to ESD, as shown in [12]. From the SLD family, we select SLD-Max, configured with an aggressive hyper-parameter setting (Hyp-Max) for inappropriate concept unlearning. It is worth noting that not all unlearned DMs are developed to address concept, style, and object unlearning tasks simultaneously. Therefore, we assess their robustness solely within the specific unlearning scenarios that they were originally designed for. By default, the victim unlearned DMs in our study are built upon Stable Diffusion (SD) v1.4. For a summary of the unlearned DMs and their corresponding unlearning tasks, please refer to Table 1.

Text Prompt Setup. In text-to-image generation, various inputs such as text prompts, random seed values, and guidance scales can be altered to generate diverse images [7]. Hence, we assess the robustness of unlearned DMs using their original prompt, random seed, and guidance scale configurations for each unlearning instance. This ensures that these victim unlearned models, without (subtle) prompt perturbations, can effectively prevent the generation of unwanted original prompt-driven images. To assess victim models' robustness in concept unlearning, we utilize the original text prompts sourced from the inappropriate image prompt (I2P) dataset [10]. This dataset targets image generation with harmful content, including nudity, violence, and illegal content. For style unlearning, we focus on artistic style and employ artistic painting prompts following the setup in [12]. In the context of object unlearning, we employ GPT-4

[64] to generate 50 text prompts for each object class in Imagenette [65]. Following a sanity check, these text prompts are confirmed to enable the vanilla SD (stable diffusion) to generate images with Imagenette objects.

Attack Setup. To address the optimization problem associated with our proposed UnlearnDiffAtk, we set the prepended prompt perturbations by N tokens, where N=5 is for concept unlearning, while N=3 is for style and object unlearning. These choices are based on the observation that prompts for concept unlearning tend to be considerably longer. To optimize the adversarial perturbations, we sample 50 diffusion time steps and perform PGD running for 40 iterations with a learning rate of 0.01 at each step. Prior to projection onto the discrete text space, we utilize the AdamW optimizer [66].

**Evaluation Metrics.** Before conducting the robustness evaluation, we have confirmed that the victim unlearned DMs are unable to generate images corresponding to an inappropriate prompt c, while their vanilla versions can successfully generate images related to c. Subsequently, to gauge the effectiveness of the proposed adversarial prompt (UnlearnDiffAtk) in steering an unlearned DM to generate undesired images (related to the original prompt c), we utilize post-generation image classifiers to determine whether the unlearning protections have been bypassed. For concept unlearning, we employ NudeNet [67] to detect nudity in the images and the Q16 classifier [68] to detect other harmful subjects. For style unlearning, we finetune an ImageNet-pretrained ViT-base model [69] on WikiArt dataset [70] to obtain a 129-class style classifier. For object unlearning, we employ the ImageNet-pretrained ResNet-50 for generated image classification. When facing inappropriate test prompts, we will dissect the attack success rate (ASR) into two categories: (1) the pre-attack success rate (pre-ASR), and (2) the post-attack success rate (post-ASR). The effectiveness of our proposed attack will be quantified by post-ASR as it measures the number of successfully bypassed unlearning safeguards using adversarial perturbations. However, a higher pre-ASR also signifies the inherent robustness of an

Table 2. Performance of various attack methods against unlearned DMs in concept unlearning, measured by attack success rate (ASR) and computation time in minutes (mins). 'No Attack' uses original prompts from I2P. 'P4D' [31] and UnlearnDiffAtk (ours) are optimization-based attack methods. 'Attack Time' represents the average computation time for generating one attack per prompt. The best attack performance (highest ASR or lowest computation time) is highlighted in bold.

I2P:			Nudity			Violence	III	Atk. Time			
Total	Prompts #:		142			756		727		per Prompt	
Unlea	rned DMs:	ESD	FMN	SLD	ESD	FMN	SLD	ESD	FMN	SLD	(mins)
Attacks:	No Attack	20.42%	88.03%	33.10%	27.12%	43.39%	22.93%	30.99%	32.83%	27.78%	-
(ASR %)	P4D	69.71%	97.89%	77.46%	80.56%	85.85%	62.43%	85.83%	<b>88.03</b> %	81.98%	34.70
	UnlearnDiffAtk	<b>76.05</b> %	97.89%	<b>82.39</b> %	80.82%	84.13%	<b>62.57</b> %	85.01%	86.66%	<b>82.81</b> %	26.29



Fig. 4. Generated images using ESD under different attacks for concept unlearning.

unlearned DM on its own. Since ASR = pre-ASR + post-ASR, we will report ASR together with pre-ASR for ease of presentation. We also remark that ASR reduces to pre-ASR when *no* adversarial attack is applied to text prompts.

#### 5.2 Experiment Results

In the following, we demonstrate from three unlearning categories (*Concept*, *Style*, *Object*) that UnlearnDiffAtk remains effective without the guidance of auxiliary models, and it improves time efficiency.

Robustness Evaluation of Unlearned DMs in Concept Unlearning. In Table 2, we present the performance of various attack methods against unlearned DMs designed to mitigate the influence of inappropriate concepts from the I2P dataset. We examine three unlearned DMs: ESD, FMN, and SLD, as shown in Table 1. Our evaluation assesses their robustness across three categories of harmful concepts: nudity, violence, and illegal activity, comprising 142, 756, and 727 inappropriate prompts, respectively. We compare the attack performance of using the proposed UnlearnDiffAtk with that of two attack baselines: 'No attack', which uses the original inappropriate prompt from I2P; and 'P4D', which corresponds to the attack proposed in [31] to solve the optimization problem (8). It is worth noting that P4D is a concurrent development while we were preparing our draft. Additionally, we compare different attack methods with respect to 'attack time' (Atk. time), given by the average computation time needed to generate one attack per prompt on a single NVIDIA RTX A6000 GPU. As we can see, the optimization-based attacks (both UnlearnDiffAtk and P4D) can effectively circumvent various types of unlearned DMs, achieving higher ASR

than 'No Attack'. Moreover, in most cases, UnlearnDiffAtk outperforms P4D although the ASR gap is not quite significant in concept learning. However, our improvement is achieved using lower computational cost than P4D, reducing runtime cost per attack instance generation by approximately 23.5%. By viewing from ASR, ESD demonstrates better robustness than other unlearned DMs, including FMN and SLD, when facing inappropriate prompts. Figure 4 displays a collection of generated images under the obtained adversarial prompts against ESD. For instance, when comparing the perturbed prompt  $P_4$  generated with UnlearnDiffAtk to the one produced with P4D, we observe that the former results in more aggressive generation. A similar pattern is observed with prompts  $P_5$  and  $P_6$ , which generate images featuring the illegal substance ('drug') and the action of 'police arrest'. More examples can be found in Fig. A2.

**Table 3.** Attack performance of various methods against unlearned DMs in Van Gogh's painting style unlearning, measured by ASR averaged over perturbing 50 Van Goghrelated prompts, and average attack time for generating one attack per prompt. The best attack performance (highest ASR or lowest attack time) is highlighted in **bold**.

Artistic	Van Gogh							Atk. Time per Prompt (mins)		
Unlearned DMs:		E	SD	FN	IN	AC UCE				
		Top-1	Top-3	Top-1	Top-3	Top-1	Top-3	Top-1	Top-3	
Attacks: (ASR %)	No Attack	2.00%	16.00%	10.00%	32.00%	12.00%	52.00%	62.00%	78.00%	-
	P4D	30.00%	78.00%	54.00%	90.00%	68.00%	94.00%	98.00%	100.00%	50.79
	UnlearnDiffAtk	<b>32.00</b> %	76.00%	56.00%	90.00%	77.00%	92.00%	94.00%	100.00%	38.87

Robustness Evaluation of Unlearned DMs in Style Unlearning. In Table 3, we present the attack performance against unlearned DMs, specifically targeting the removal of the 'Van Gogh's painting style' influence in image generation. This style of unlearning has also been studied by other unlearning methods, as shown in Table 1. Unlike concept unlearning, our evaluation of ASR considers two types: 'Top-1 ASR' and 'Top-3 ASR'. These metrics depend on whether  $_{
m the}$ generated

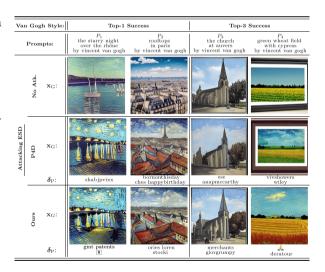


Fig. 5. Generated images using ESD under different attacks for style unlearning.

image ranks as the top-1 prediction or within the top-3 predictions regarding Van Gogh's painting style when assessed by the post-generation image classifier. This is motivated by our observation that relying solely on the top-1 prediction

might be overly restrictive when assessing the relevance to Van Gogh's painting style; See Fig. 5. Moreover, consistent with [12], we employ 50 prompts for image generation with the Van Gogh style and utilize them to assess the robustness of unlearned DMs. Similar to Table 2, we compare our proposed UnlearnDiffAtk with 'no attack' and P4D on four unlearned DMs: ESD, FMN, AC, and UCE. As we can see, UnlearnDiffAtk continues to prove its effectiveness and efficiency as an attack method to bypass the unlearned DMs, enabling the generation of images with the Van Gogh's painting style. Among the unlearned DMs, ESD exhibits the highest unlearning robustness when considering Top-1 ASR. Nevertheless, Top-3 ASR still maintains a performance level exceeding 80% when employing UnlearnDiffAtk, and is sufficient to indicate the generation of images with the Van Gogh's painting style, as illustrated in Fig. 5. We observe that in the absence of an attack against ESD, the generated images (e.g., under  $P_4$ ) lack Van Gogh's painting style. However, UnlearnDiffAtk-enabled prompt perturbations can effortlessly bypass ESD, resulting in the generation of Van Gogh-style images. More generated images can be found in Fig. A3.

**Table 4.** Attack performance of various methods against unlearned DMs in object unlearning, measured by ASR averaged over perturbing 50 prompts for each object class, and the average computation time for generating one attack per prompt. The best attack performance (highest ASR or lowest attack time) is highlighted in **bold**.

Object Classes:		Church Parachute		Tench		Garbage Truck		Atk. Time per Prompt (mins)		
Unlearned	l DMs:	ESD	FMN	ESD	FMN	ESD	FMN	ESD	FMN	
Attacks: (ASR %)	No Attack	14%	52%	4%	46%	2%	42%	2%	40%	-
	P4D	56%	98%	48%	100%	28%	96%	20%	98%	43.65
	UnlearnDiffAtk	60%	96%	54%	100%	36%	100%	24%	98%	31.32

Robustness Evaluation of Unlearned  $\mathrm{DMs}$ in Object Unlearning. In Table 4, we present the results showcasing the performance of different attacks concerning object unlearning. We regard ESD and FMN as the victim models, which erase one of the chosen four object classes from Imagenette [65]. These specific classes were selected due to their ease of differentiation, allowing us to assess the effectiveness of the attacks.

Object	Classes:	Church	Parachute	Tench	Garbage Truck
Pro	mpts:	P <sub>1</sub> Church surrounded by autumn foliage.	$P_2$ Parachute in a desert landscape.	P <sub>3</sub> Tench in a fish market.	P <sub>4</sub> Garbage truck during winter.
No Atk	$\mathbf{x}_{\mathrm{G}}$ :		• • •		
Attacking ESD P4D	$\mathbf{x}_{\mathrm{G}}$ : $\boldsymbol{\delta}_{\mathrm{P}}$ :	blanc sheep ges	bersersings confrontation	qe wicked atlanta	matteo yelling promot
Ours		hoengineerhain	wrinkles staining modest	itf @ mixed	trunks personnel waxin

Fig. 6. Generated images using ESD under different attacks for object unlearning.

Given an image class, we apply each attack method to 50 prompts generated using ChatGPT that pertain to this class. Similar to concept and style unlearning, we compare the ASR and the attack generation time of UnlearnDiffAtk with 'No Attack' and P4D. As we can see, UnlearnDiffAtk consistently achieves a higher ASR than P4D across various unlearning objects and victim models while requiring less computational resources. Furthermore, ESD demonstrates better robustness against prompt perturbations than FMN in the context of object unlearning. Figure 6 displays generation examples under the obtained adversarial prompts against ESD. We note that the objects (such as 'Parachute' in  $P_2$  and 'Garbage Truck' in  $P_4$ ) can be re-generated under UnlearnDiffAtk-perturbed prompts, as compared to P4D and No Attack. More results can be found in Fig. A4.

Attack Using Different Target Image Sources. As discussed in Remark 1 of Sect. 4, our proposed UnlearnDiffAtk benefits from its sole reliance on a target image  $\mathbf{x}_{\mathrm{tgt}}$ , without requiring an auxiliary vanilla DM during attack generation. In our prior

Table 5. ASR of UnlearnDiffAtk when attacking ESD (based on SD v1.4) using target images generated from either SD v1.4 or SD v2.1.

UnlearnDiffAtkvs. ES	Nudity	Van Gogh		Church	
			Top-1	Top-3	
DM of Target	SD v1.4	76.05%	32.00%	76.00%	60.00%
Image Generation	SD v2.1	73.94%	34.00%	82.00%	60.00%

experiments, we explored this setting with  $\mathbf{x}_{tgt}$  generated using SD v1.4, the same SD version used by unlearned DMs. Table 5 shows the ASR achieved when utilizing UnlearnDiffAtk against the ESD model (built upon SD v1.4), given that the target image  $\mathbf{x}_{tgt}$  is generated using different versions of SD, v1.4 and v2.1, respectively. We observe that UnlearnDiffAtk maintains a consistent ASR performance, even when there's a discrepancy between the target image source (acquired by SD v2.1) and the victim model, ESD built upon SD v1.4.

Other Ablation Studies. In Appendix B, we demonstrate more ablation studies. This includes (1) the resilience of attack performance against the adversarial prompt location and length (Table A1 and Table A2), (2) the attack transferability across different SD models (Table A3), and (3) attack effectiveness compared to 'random' attacks (Table A4).

### 6 Conclusions

The evolution of DMs (diffusion models) in generating intricate images underscores both their potential and their inherent risks. While these models present significant advancements in the realm of digital imagery, the capacity for generating unsafe content cannot be understated. Our research sheds light on the vulnerabilities of current safety-driven unlearned DMs when confronted with adversarial prompts, even when these prompts involve subtle text perturbations. Notably, we develop the UnlearnDiffAtk method, which not only simplifies the generation of adversarial prompts against DMs (without the need of auxiliary models) but also offers an innovative perspective on utilizing DMs'

classification capabilities. We also conduct a comprehensive set of experiments to benchmark the robustness of state-of-the-art unlearned DMs across multiple unlearning tasks. Our research emphasizes the need for more resilient and trustworthy systems in conditional diffusion-based image generation systems.

**Acknowledgement.** Y. Zhang, J. Jia, Y. Zhang, C. Fan, J. Liu, and S. Liu were supported by the National Science Foundation (NSF) Robust Intelligence (RI) Core Program Award IIS-2207052, the NSF CPS Award CNS-2235231, the Cisco Faculty Research Award, and the ARO Award W911NF2310343.

### References

- Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Adv. Neural. Inf. Process. Syst. 33, 6840–6851 (2020)
- Song, Y., Ermon, S.: Generative modeling by estimating gradients of the data distribution. In: Advances in Neural Information Processing Systems, vol. 32 (2019)
- 3. Song, Y., Sohl-Dickstein, J., Kingma, D.P., Kumar, A., Ermon, S., Poole, B.: Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456 (2020)
- Dhariwal, P., Nichol, A.: Diffusion models beat GANs on image synthesis. In: Advance in Neural Information Processing System, vol. 34, pp. 8780–8794 (2021)
- Nichol, A.Q., Dhariwal, P.: Improved denoising diffusion probabilistic models. In: International Conference on Machine Learning, pp. 8162–8171. PMLR (2021)
- Watson, D., Chan, W., Ho, J., Norouzi, M.: Learning fast samplers for diffusion models by differentiating through sample quality. In: International Conference on Learning Representations (2021)
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10684–10695 (2022)
- 8. Croitoru, F.A., Hondru, V., Ionescu, R.T., Shah, M.: Diffusion models in vision: a survey. IEEE Trans. Pattern Anal. Mach. Intell. (2023)
- 9. Rando, J., Paleka, D., Lindner, D., Heim, L., Tramèr, F.: Red-teaming the stable diffusion safety filter. arXiv preprint arXiv:2210.04610 (2022)
- Schramowski, P., Brack, M., Deiseroth, B., Kersting, K.: Safe latent diffusion: mitigating inappropriate degeneration in diffusion models (2023)
- 11. Nichol, A., et al.: Glide: towards photorealistic image generation and editing with text-guided diffusion models. arXiv preprint arXiv:2112.10741 (2021)
- Gandikota, R., Materzynska, J., Fiotto-Kaufman, J., Bau, D.: Erasing concepts from diffusion models. arXiv preprint arXiv:2303.07345 (2023)
- Brack, M., Friedrich, F., Schramowski, P., Kersting, K.: Mitigating inappropriateness in image generation: Can there be value in reflecting the world's ugliness? arXiv preprint arXiv:2305.18398 (2023)
- 14. Yang, Y., Hui, B., Yuan, H., Gong, N., Cao, Y.: Sneakyprompt: Evaluating robustness of text-to-image generative models' safety filters. arXiv preprint arXiv:2305.12082 (2023)
- Zhang, E., Wang, K., Xu, X., Wang, Z., Shi, H.: Forget-me-not: learning to forget in text-to-image diffusion models. arXiv preprint arXiv:2303.17591 (2023)
- Kumari, N., Zhang, B., Wang, S.Y., Shechtman, E., Zhang, R., Zhu, J.Y.: Ablating concepts in text-to-image diffusion models (2023)

- 17. Gandikota, R., Orgad, H., Belinkov, Y., Materzyńska, J., Bau, D.: Unified concept editing in diffusion models. arXiv preprint arXiv:2308.14761 (2023)
- 18. Nguyen, T.T., Huynh, T.T., Nguyen, P.L., Liew, A.W.C., Yin, H., Nguyen, Q.V.H.: A survey of machine unlearning. arXiv preprint arXiv:2209.02299 (2022)
- Shaik, T., Tao, X., Xie, H., Li, L., Zhu, X., Li, Q.: Exploring the landscape of machine unlearning: a survey and taxonomy. arXiv preprint arXiv:2305.06360 (2023)
- 20. Cao, Y., Yang, J.: Towards making systems forget with machine unlearning. In: 2015 IEEE Symposium on Security and Privacy, pp. 463–480. IEEE (2015)
- 21. Thudi, A., Deza, G., Chandrasekaran, V., Papernot, N.: Unrolling SGD: understanding factors influencing machine unlearning. In: 2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P), pp. 303–319. IEEE (2022)
- Xu, H., Zhu, T., Zhang, L., Zhou, W., Yu, P.S.: Machine unlearning: a survey.
   ACM Comput. Surv. 56(1), 1–36 (2023)
- 23. Jia, J., et al.: Model sparsity can simplify machine unlearning (2023)
- 24. Zhang, Y., et al.: Unlearncanvas: a stylized image dataset to benchmark machine unlearning for diffusion models. arXiv preprint arXiv:2402.11846 (2024)
- Jia, J., et al.: Soul: unlocking the power of second-order optimization for LLM unlearning. arXiv preprint arXiv:2404.18239 (2024)
- Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572 (2014)
- 27. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: 2017 IEEE Symposium on Security and Privacy (SP), pp. 39–57. IEEE (2017)
- Maus, N., Chao, P., Wong, E., Gardner, J.R.: Black box adversarial prompting for foundation models. In: The Second Workshop on New Frontiers in Adversarial Machine Learning (2023)
- Zhuang, H., Zhang, Y., Liu, S.: A pilot study of query-free adversarial attack against stable diffusion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2384–2391 (2023)
- 30. Wen, Y., Jain, N., Kirchenbauer, J., Goldblum, M., Geiping, J., Goldstein, T.: Hard prompts made easy: gradient-based discrete optimization for prompt tuning and discovery. arXiv preprint arXiv:2302.03668 (2023)
- 31. Chin, Z.Y., Jiang, C.M., Huang, C.C., Chen, P.Y., Chiu, W.C.: Prompting4debugging: red-teaming text-to-image diffusion models by finding problematic prompts. arXiv preprint arXiv:2309.06135 (2023)
- 32. Birhane, A., Prabhu, V.U., Kahembwe, E.: Multimodal datasets: misogyny, pornography, and malignant stereotypes. arXiv preprint arXiv:2110.01963 (2021)
- Somepalli, G., Singla, V., Goldblum, M., Geiping, J., Goldstein, T.: Diffusion art or digital forgery? investigating data replication in diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6048–6058 (2023)
- 34. Schuhmann, C., et al.: Laion-5b: an open large-scale dataset for training next generation image-text models. In: Advance in Neural Information Processing System, vol. 35, pp. 25278–25294 (2022)
- 35. Kumari, N., Zhang, B., Wang, S.Y., Shechtman, E., Zhang, R., Zhu, J.Y.: Ablating concepts in text-to-image diffusion models. In: ICCV (2023)
- 36. Heng, A., Soh, H.: Selective amnesia: a continual learning approach to forgetting in deep generative models. arXiv preprint arXiv:2305.10120 (2023)
- 37. Ni, Z., Wei, L., Li, J., Tang, S., Zhuang, Y., Tian, Q.: Degeneration-tuning: using scrambled grid shield unwanted concepts from stable diffusion. arXiv preprint arXiv:2308.02552 (2023)

- 38. Zhang, Y., et al.: Defensive unlearning with adversarial training for robust concept erasure in diffusion models. arXiv preprint arXiv:2405.15234 (2024)
- 39. Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z.B., Swami, A.: The limitations of deep learning in adversarial settings. In: Security and Privacy (EuroS&P), 2016 IEEE European Symposium on, pp. 372–387. IEEE (2016)
- Brown, T.B., Mané, D., Roy, A., Abadi, M., Gilmer, J.: Adversarial patch. arXiv preprint arXiv:1712.09665 (2017)
- 41. Li, J., Schmidt, F., Kolter, Z.: Adversarial camera stickers: a physical camerabased attack on deep learning systems. In: International Conference on Machine Learning, pp. 3896–3904 (2019)
- 42. Xu, K., et al., Lin, X.: Structured adversarial attack: towards general implementation and better interpretability. In: ICLR (2019)
- 43. Yuan, Z., Zhang, J., Jia, Y., Tan, C., Xue, T., Shan, S.: Meta gradient adversarial attack. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 7748–7757 (2021)
- 44. Zhang, Y., Yao, Y., Jia, J., Yi, J., Hong, M., Chang, S., Liu, S.: How to robustify black-box ml models? a zeroth-order optimization perspective. arXiv preprint arXiv:2203.14195 (2022)
- Chen, A., et al.: Deepzero: scaling up zeroth-order optimization for deep model training. arXiv preprint arXiv:2310.02025 (2023)
- 46. Gong, Y., et al.: Reverse engineering of imperceptible adversarial image perturbations. arXiv preprint arXiv:2203.14145 (2022)
- 47. Qiu, S., Liu, Q., Zhou, S., Huang, W.: Adversarial attack and defense technologies in natural language processing: a survey. Neurocomputing 492, 278–307 (2022)
- 48. Eger, S., Benz, Y.: From hero to zéroe: a benchmark of low-level adversarial attacks. In: Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing, pp. 786–803 (2020)
- 49. Liu, A., et al.: Character-level white-box adversarial attacks against transformers via attachable subwords substitution. arXiv preprint arXiv:2210.17004 (2022)
- Hou, B., et al.: TextGrad: advancing robustness evaluation in NLP by gradientdriven optimization. arXiv preprint arXiv:2212.09254 (2022)
- 51. Li, J., Ji, S., Du, T., Li, B., Wang, T.: TextBugger: generating adversarial text against real-world applications. arXiv preprint arXiv:1812.05271 (2018)
- Alzantot, M., Sharma, Y., Elgohary, A., Ho, B.J., Srivastava, M., Chang, K.W.: Generating natural language adversarial examples. arXiv preprint arXiv:1804.07998 (2018)
- 53. Jin, D., Jin, Z., Zhou, J.T., Szolovits, P.: Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 34, pp. 8018–8025 (2020)
- 54. Garg, S., Ramakrishnan, G.: Bae: bert-based adversarial examples for text classification. arXiv preprint arXiv:2004.01970 (2020)
- Pham, M., Marshall, K.O., Cohen, N., Mittal, G., Hegde, C.: Circumventing concept erasure methods for text-to-image generative models. In: The Twelfth International Conference on Learning Representations (2023)
- Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or, D.: An image is worth one word: personalizing text-to-image generation using textual inversion. arXiv preprint arXiv:2208.01618 (2022)
- 57. Cao, H., et al.: A survey on generative diffusion model. arXiv preprint arXiv:2209.02646 (2022)

- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083 (2017)
- Croce, F., Hein, M.: Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In: ICML (2020)
- Chen, H., et al.: Robust classification via a single diffusion model. arXiv preprint arXiv:2305.15241 (2023)
- Li, A.C., Prabhudesai, M., Duggal, S., Brown, E., Pathak, D.: Your diffusion model is secretly a zero-shot classifier. arXiv preprint arXiv:2303.16203 (2023)
- Iusem, A.N.: On the convergence properties of the projected gradient method for convex optimization. Comput. Appl. Math. 22, 37–52 (2003)
- 63. Parikh, N., Boyd, S., et al.: Proximal algorithms. Found. Trends® Optim. 1(3), 127–239 (2014)
- 64. OpenAI: Gpt-4 technical report. ArXiv abs/2303.08774 (2023). https://api.semanticscholar.org/CorpusID:257532815
- Shleifer, S., Prokop, E.: Using small proxy datasets to accelerate hyperparameter search. arXiv preprint arXiv:1906.04887 (2019)
- 66. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017)
- 67. Bedapudi, P.: Nudenet: neural nets for nudity classification, detection and selective censoring (2019)
- 68. Schramowski, P., Tauchmann, C., Kersting, K.: 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, pp. 1350–1361 (2022)
- 69. Wu, B., et al.: Visual transformers: token-based image representation and processing for computer vision (2020)
- 70. Saleh, B., Elgammal, A.: Large-scale classification of fine-art paintings: learning the right metric on the right feature. arXiv preprint arXiv:1505.00855 (2015)