# LoRAShield: Data-Free Editing Alignment for Secure Personalized LoRA Sharing

Jiahao Chen, Junhao Li, Yiming Wang, Zhe Ma, Yi Jiang, Chunyi Zhou,
Qingming Li, Tianyu Du, Shouling Ji
College of Computer Science and Technology, Zhejiang University

## ABSTRACT

The proliferation of Low-Rank Adaptation (LoRA) models has democratized personalized text-to-image generation, enabling users to share lightweight models (e.g., personal portraits) on platforms like Civitai and Liblib. However, this "share-and-play" ecosystem introduces critical risks: benign LoRAs can be weaponized by adversaries to generate harmful content (e.g., political, defamatory imagery), undermining creator rights and platform safety. Existing defenses like concept-erasure methods focus on full diffusion models (DMs), neglecting LoRA's unique role as a modular adapter and its vulnerability to adversarial prompt engineering. To bridge this gap, we propose LoRAShield, the first data-free editing framework for securing LoRA models against misuse. Our platform-driven approach dynamically edits and realigns LoRA's weight subspace via adversarial optimization and semantic augmentation. Experimental results demonstrate that LoRAShield achieves remarkable effectiveness, efficiency, and robustness in blocking malicious generations without sacrificing the functionality of the benign task. By shifting the defense to platforms, LoRAShield enables secure, scalable sharing of personalized models, a critical step toward trustworthy generative ecosystems. **Warnings**: *This paper contains sexually and bloody explicit imagery that some readers may find disturbing, distressing, and/or offensive. To mitigate the offensiveness to readers, we showcase some of the explicit images with black masks.*

## CCS CONCEPTS

• **Security and privacy** → *Digital rights management*; • **Do Not Use This Code** → **Generate the Correct Terms for Your Paper**; *Generate the Correct Terms for Your Paper*; Generate the Correct Terms for Your Paper; Generate the Correct Terms for Your Paper.

## KEYWORDS

diffusion model, text-to-image, adversarial attack, concept erasure

(a) Misuse example 1    (b) Misuse example 2

**Figure 1: Generated images using LoRAs ([18] and [35]) on Civitai. Note that these images will never be posted online.**

## 1 INTRODUCTION

The rapid development of text-to-image (T2I) technology [15, 38], particularly diffusion models (DMs) [41], has dramatically improved personalized creativity. However, this advancement has simultaneously exacerbated problems like bias [16, 32], copyright infringement [44], and the generation of not-safe-for-work (NSFW) content [51, 52]. Malicious users of T2I services may exploit these models to create misleading images that harm individuals' interests or reputations or even provoke social panic [39]. For instance, a finance worker at a multinational firm was tricked into paying out $25 million to fraudsters using deepfake technology to pose as the company's chief financial officer in a video conference call [4].

To mitigate such risks, numerous studies [20, 49, 58] have explored the elimination of specific concepts from DMs to prevent the generation of harmful content. However, existing research focuses mainly on the DMs themselves [57], overlooking significant vulnerabilities in their plugins. In particular, Low-Rank Adaptation (LoRA) models [23] enable parameter-efficient adaptation of DMs through low-rank matrix decomposition, allowing users to customize models with minimal computational overhead. The widespread adoption of LoRA, due to its lightweight characteristics, enables individuals to efficiently build personalized models from their local datasets, such as personal portraits or artworks, utilizing DMs as base models. Users often share and publish their LoRA models, instead of the entire DM, on platforms like Civitai [8], LiblibAI [29], and Hugging Face [24]. Any other user can download these shared LoRAs and incorporate them into the generation of fancier images.

Unfortunately, this share-and-play characteristic also introduces significant risks. Malicious users can misuse others' uploaded LoRAs to generate harmful content. For example, our tests in Fig. 1 using Civitai's online generation services reveal that despite the platform's claims of preventing the generation of NSFW content [9], benign LoRAs uploaded by genuine users can still be exploited to produce content that can significantly disrupt the interest and reputation of LoRA owners. Thus, malicious actors can download benign LoRAs and generate harmful images, negatively impacting the

original creators' interests or reputation. This situation indicates a critical gap in effective oversight by both academia and industry regarding the abuse of LoRAs.

To address this risk, we re-examined the threat models (Section 3.3) associated with LoRA misuse in current model marketplaces and wondered if it is possible to develop a defense to mitigate the gap above. Given the high volume of LoRA models uploaded to platforms [47], the defense method must be lightweight and efficient, avoiding negative impacts on the quality of legitimate image generation. Additionally, we find that even though previous studies [49, 57] have indeed proposed methods for concept erasure in DMs, these typically require substantial computational resources and extra images. Moreover, current concept erasure techniques frequently exhibit instability, incomplete removal of target concepts, limited generalization, and robustness restricted to specific prompts, which will be discussed in Section 2.2 as well.

Further, we propose a platform-driven and data-free LoRA editing alignment method (LoRAShield) that allows users to specify undesired concepts, thus preventing their LoRAs from being exploited for malicious purposes. Specifically, users upload LoRAs to a platform that offers dynamic editing services based on user-defined constraints. By editing and realigning the LoRA's weight subspace, LoRAShield restricts the model's ability to produce certain undesired content while preserving the quality of its intended generation tasks for legitimate uses. In summary, we make the following contributions.

- We reveal the risk of the current personalization model on model-sharing platforms, raising the alarm about harmful image generation for LoRA owners.
- We make the first attempt to develop a data-free editing method to secure the sharing of personalized LoRA with adversarial optimization and semantic augmentation.
- Extensive experiments demonstrate the effectiveness, robustness, and efficiency of LoRAShield.

## 2 BACKGROUND & RELATED WORK

### 2.1 Model-Sharing Platforms

The proliferation of model-sharing platforms such as Civitai [8], LiblibAI [29], and Hugging Face [24] has democratized access to generative model tools, fostering a "share and play" culture where users freely distribute and utilize pre-trained models. However, this accessibility introduces critical risks. While this facilitates innovation, it also creates vulnerabilities: benign users who upload models trained on private data face potential misuse by adversaries. Such risks are exacerbated by the lack of robust safeguards on many platforms; for instance, Civitai's reliance on an "asymmetric multisided marketplace" [21] prioritizes model diversity over rigorous vetting.

### 2.2 Concept Erasing in DMs

Large-scale DMs are shadowed by their propensity to produce ethically risky or legally contentious outputs, such as sexually explicit imagery, culturally sensitive content, or artistic styles protected by copyright. A growing number of research studies [17, 19, 20, 27, 53, 54, 57] have proposed leveraging unlearning techniques to remove or suppress specific concepts via unlearning in generative models to prevent misuse. The traditional machine unlearning

strategy [2] aims to modify a model to make it generalize without memorization to tackle privacy and copyright concerns. Although early approaches [19, 26] focused on fine-tuning all DMs parameters to forget a target concept. For example, the Erased Stable Diffusion (ESD) by Gandikota et al.[19, 26, 53] fine-tunes a DM with a negative guidance teacher to align the given visual concept with " " condition. Similarly, Zhang et al.[53] proposed a "Forget-Me-Not" (FMN) technique that aligns the model output for the target concept with that of a benign "anchor" concept, thus preventing the model from generating the target concept under its text condition, allowing the model to forget without retraining from scratch while preserving closely related concepts. However, since these methods tend to struggle with sizeable computational costs [20, 57], Unified Concept Editing (UCE) [20] proposed to edit DMs without training to improve the efficacy and scalability of the previous work. Wu et al.[49] formulated diffusion unlearning as a constraint optimization and achieved it by deviating the learnable generative process from the ground-truth denoising procedure, while the assumption requiring all training data is unrealistic. Later, Lyu et al. [33] argued that most of the existing methods failed to preserve the generation of untargeted concepts and proposed an erasing framework (SPM) via one-dimensional adapters to erase multiple concepts from DMs at once via versatile erasing. Differently, SafeGen [27] eliminated unsafe visual representations from the model regardless of text input, leading to resistance to adversarial prompts since unsafe visual representations are obstructed from within.

### 2.3 Adversarial Attacks against DMs

The growing adoption of defensive strategies in DMs has catalyzed parallel research into system robustness and security, revealing critical vulnerabilities through adversarial attacks. Recent work demonstrates that DMs remain susceptible to generating harmful content even after deploying safety measures, particularly when adversaries craft inputs to circumvent safeguards. For example, Yang et al. [52] crafts semantically preserved yet misleading prompts on surrogate models, enabling transferable attacks that deceive prompt filters. Instead, MMA-Diffusion [51] conducted transferable attacks by crafting adversarial prompts on surrogate models that can deceive the prompt filter while remaining semantically similar to the target concept. Zhang et al. [56] further exposed the fragility of concept-erasure techniques by proposing UnlearnDiffAtk, an adversarial prompt generation method that resurrects supposedly "erased" concepts in fine-tuned DMs. Complementing these efforts, automated red-teaming tools such as Ring-A-Bell[45] and Prompting4Debugging [6] systematically stress-test DM safety mechanisms. These papers reveal that even DMs subjected to rigorous fine-tuning or filtering remain vulnerable to synonym substitutions, highlighting that superficial compliance with safety benchmarks does not guarantee robustness against adaptive adversaries.

### 2.4 Remarks

Overall, the existing literature indicates the following points: (1) all of the concept-erasing methods target DMs instead of LoRAs, which are frequently used for customization of individuals; (2) recent evaluation [57] also reveals that most of these attacks exhibit limited generalization to in-domain prompts (innocent and not relevant to

**Table 1: Comparison with existing erasing methods. ✔/✗ illustrates whether the method can achieve the corresponding property.**

| Erasing Methods | Target | Effectiveness | | | Efficiency | |
|---|---|---|---|---|---|---|
| | | Generalization | Robustness | Data-Free | Time (s) | Memory (GB) |
| ESD [19] | DMs | ✔ | ✗ | ✔ | $\approx 6100$ | 17.8 |
| FMN [53] | DMs | ✗ | ✗ | ✗ | $\approx 350$ | 17.9 |
| UCE [20] | DMs | ✗ | ✗ | ✔ | $\approx 430$ | 5.1 |
| EraseDiff [49] | DMs | ✔ | ✗ | ✗ | $\approx 1500$ | 27.8 |
| SPM [33] | DMs | ✔ | ✗ | ✗ | $\approx 29700$ | 6.9 |
| SafeGen [27] | DMs | ✔ | ✔ | ✗ | – | – |
| Ours (LoRAShield) | LoRA | ✔ | ✔ | ✔ | $\approx 14$ | 0.23 |

the erased concept); (3) the majority of these methods fail to defense against current attacks [51, 52, 56]; (4) balancing completeness of erasure with minimal data, time and computation resources is still challenging. In comparison, as shown in Tab. 1, LoRAShield addresses these gaps and secures the modular components most vulnerable to misuse in platforms.

## 3 PRELIMINARIES

### 3.1 Text-to-Image Diffusion Models

T2I diffusion models [15] can synthesize high-fidelity images $x \in \mathcal{X}$ by iteratively denoising latent representations $z_t \in \mathcal{Z}$ conditioned on textual inputs $c$. The core framework typically comprises four components: an image encoder $\mathcal{E} : \mathcal{X} \rightarrow \mathcal{Z}$ and decoder $\mathcal{D} : \mathcal{Z} \rightarrow \mathcal{X}$ of the autoencoder, a U-Net diffusion model $\epsilon_\theta$, and a text encoder $\mathcal{T}$. The encoder first compresses images into a lower-dimensional latent space and reconstructs them via a decoder $\mathcal{D}$, i.e., $z = \mathcal{E}(x)$ and $\hat{x} = \mathcal{D}(z)$. While the U-Net then operates in this latent space, iteratively refining the noisy representations $z_t$ through a reverse diffusion process:

$$z_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( z_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \cdot \epsilon_\theta(z_t, t, c) \right) + \sigma_t \cdot \epsilon, \qquad (1)$$

where timestep $t \in \{1, 2, \ldots, T\}$, $\epsilon \sim \mathcal{N}(0, I)$, $\epsilon_\theta$ predicts the noise, $c = \mathcal{T}(C)$ denotes text embeddings given the prompt set $C$, and $\alpha_t, \bar{\alpha}_t, \sigma_t$ are scheduler hyperparameters. Specifically, the U-Net integrates textual guidance via cross-attention layers, which dynamically align image features with target concepts. Cross-attention (CA) layer computes:

$$\text{Atten}(Q, K, V) = \text{softmax}\left( \frac{QK^T}{\sqrt{d_k}} \right) V, \qquad (2)$$

where queries $Q$ derive from image tokens, and keys $K$ and values $V$ derive from $c$. This mechanism enables fine-grained control over concept representation, as the model prioritizes text-relevant features during denoising.

### 3.2 Low-Rank Adaptation of Diffusion Models

The LoRA techniques [23] enable efficient and effective adaptation and customization of T2I diffusion models via decomposing the weight update matrix $\Delta W \in \mathbb{R}^{m \times n}$ when fine-tuning the original weight $W_0$:

$$\Delta W = B \times A \qquad (3)$$

where $B \in \mathbb{R}^{m \times r}$, $A \in \mathbb{R}^{r \times n}$ and $r \ll \min(m, n)$ stands for the rank that constrains the update to specified dimension space, significantly reducing the trainable parameters of $W_0$. During the inference stage, $W_0$ is updated with $W_0 + \alpha \Delta W$, where $\alpha$ is a scaling factor that control the strength of the adaptation.

The rise of platforms like Civitai [8] has democratized access to personalized generative models. Unlike full diffusion models (e.g., SD1.5 [12] and its extended versions like DreamShaper [10] and Realistic Vision [11]), which are prohibitively large (often more than 5GB) and computationally expensive to share, LoRA models compress adaptation into lightweight matrices (typically about 100MB). Users train LoRAs on local data (e.g., personal artwork or portraits) and upload them with metadata specifying compatible base models. There are much more LoRAs but less base models on Civitai [8]. This modular approach allows others to: (1) Efficiently combine concepts: Load multiple LoRAs (e.g., "cyberpunk style" + "celestial lighting") into a single base model, scaling their influence via $\alpha$; (2) Preserve reproducibility: Base models act as standardized backbones, ensuring consistent generation across users.

### 3.3 Threat Model

We consider a scenario involving three parties as depicted in Fig. 2: a LoRA owner (victim) who wants to share her/his customized LoRA trained with private data, a trusted model-sharing platform (defender) [9] and a LoRA misuse infringer (attacker).

*3.3.1 LoRA owner's Goal & Capabilities .* The LoRA owner aims to release his/her customized LoRA trained with private images (personal portrait and promoted products) to the public for benign purposes (e.g., personalized photography or advertisement). However, she/he does not want the models to be used for malicious purposes that may affect her/his interests and reputation, e.g., exploiting released LoRA to generate fake images to fabricate rumors about the LoRA owner. Since the owner has no control over the LoRA downloaded by others, to avoid infringements, the owner could request the platform for active editing services to erase the unwanted concept before releasing the LoRA.

*3.3.2 Adversary's Goal & Capabilities.* An attacker on this platform could be anyone who downloads a publicly available LoRA model and attempts to use it for harmful content generation that may violate the platform's content policies [9, 30]. The generated content [9, 30] could include violent imagery, sexual or pornographic
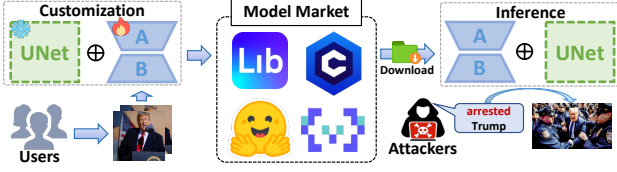
**Figure 2: Overview of the threat model.**

scenes, deepfakes of real individuals, propaganda, or other disallowed material, leveraging the specific capabilities of the LoRA. Since the attacker can download any LoRAs from the platform, the attacker (1) has complete access and control to the protected LoRA, (2) has complete control over the generation process, but (3) has no access to the training data or unprotected LoRA. However, she/he also wants to ensure the high quality and semantic alignment of the generated images.

*3.3.3 Adversary's Strategies.* We assume that the attacker is motivated and technically savvy, willing to launch adaptive attacks to bypass the protection. Note that even if most DMs like SD have safety checkers but offline attackers can turn them off or find novel prompts to evade them. (**Attack I**) One likely strategy is prompt engineering, trying various descriptive prompts (or even gibberish tokens) to trick the model into revealing the protected concept (if the LoRA was meant to be restricted). This aligns with the red-teaming approaches in recent research [52, 56], where attackers find that even "safe" models can produce unsafe outputs with cleverly crafted prompts. (**Attack II**) Another powerful strategy is to combine multiple LoRAs since these adapters are often composable, and an adversary might chain a benign LoRA with another one that introduces unsafe elements. For example, they might apply an "NSFW unlock" [13] or violence-enhanced LoRA alongside the protected LoRA, hoping to restore any pruned capabilities or amplify forbidden content. By blending LoRAs, the attacker tries to circumvent a single LoRA's built-in safety, even if our defense cripples the LoRA's ability to generate certain content on its own, the attacker could supplement it with another malicious adapter that fills in the gaps.

*3.3.4 Defender's Goal & Capabilities.* The defender's (platform) goal is to enable creative uses of LoRA while preventing misuse. For instance, a LoRA might add a new artistic style, and the platform wants regular users to enjoy it but disallows anyone from using it to generate hate symbols or obscene content in that style. Neither the benign users nor the defender trusts every user and thus implements protections specified by the LoRA owner. Therefore, the defender cannot expend significant resources and defense must be automated and efficient (no lengthy retraining or manual intervention for each model). This is because on a large platform like Civitai, there could be thousands of uploaded LoRAs [47] in a short period, and users expect to upload and update them frequently. Besides, the defender is not accessible to the benign training data of the LoRA for privacy reason, which means that the editing should be data-free as well.

## 4 METHODOLOGY

### 4.1 Design Overview

Under the threat model presented above, our primary goal is to provide an efficient and robust defense mechanism capable of preventing unauthorized and malicious exploitation of large-scale shared LoRAs on platforms while maintaining their original generative capabilities for legitimate users. In summary, our method addresses several critical challenges below:

**(1) Data-Free Implementation:** Since the platform hosting the LoRAs don't have access to the original data, the proposed method must operate entirely without reliance on the training dataset.

**(2) Preservation of Utility:** A practical defense should precisely remove only potentially harmful concept representations without diminishing the model's general utility for creative generation.

**(3) Lightweight and Efficiency:** Given the high frequency and large volume of LoRAs uploaded, our defense must be computationally lightweight to process numerous models quickly without introducing significant latency or computational overhead.

**(4) Robustness Against Adaptive Attacks:** Adversaries are aware of defenses and can employ strategies (Section 3.3.3) to evade. Hence, the defender must proactively anticipate such behaviors.

To overcome these challenges, in the following subsections, we design a lightweight framework that effectively balances security needs with usability, ensuring that LoRAs remain valuable tools while raising the barrier to malicious exploitation.

### 4.2 Robust Editing Alignment

As mentioned above in Section 3.1, given the conditional prompt $C$, the cross-attention layer matters a lot in controlling the image generation. This motivates us to erase the concept in the cross-attention layer to safeguard LoRAs. Initially, we obtain the text embeddings $c$ and $c_t$, with the benign and target concepts to be erased, respectively. Further, we calculate the difference between the LoRA-adjusted and original weight matrices for each CA layer. We iteratively update these matrices to minimize the discrepancy between the generated embeddings from the adjusted model and the intended safe embeddings, employing a loss function that combines mean squared error (MSE).

$$\mathcal{L}_{align}(c_t, c) = \mathbb{E}\|c_t \times (W + \alpha\Delta\hat{W}) - c \times (W + \alpha\Delta W)\|_2^2 \quad (4)$$

However, minimizing $\mathcal{L}_{align}$ without restricting the difference between the edited and original LoRA parameters disrupts the benign performance of the LoRA as well. Therefore, we add the regularization term to preserve the performance of the edited LoRA on the untargeted concept.

$$\mathcal{L}_{pre} = \|\Delta\hat{W} - \Delta W\|_2^2 \quad (5)$$

This targeted pruning ensures minimal disruption to benign concept generation while effectively disabling restricted concepts.

*4.2.1 Adversarial Optimization for Robust Generalization.* However, incorporating $\mathcal{L}_{align}$ and $\mathcal{L}_{pre}$ can only prevent the generation given the exact same word of the erased concept, but would fail to generalize to semantic space. Ideally, our goal is to achieve robust alignment by erasing the semantics of the target concept. To achieve this, we first make the following assumption.

**Assumption 1:** Let $w_t^i \in \mathcal{W}_t$ be the set of synonyms with the same concept in the text space. We assume that $\forall w_t^i, w_t^j \in \mathcal{W}_t$ and $\exists \gamma > 0$, and we have $\text{Dist}(\mathcal{T}(w_t^i), \mathcal{T}(w_t^j)) \leq \gamma$, where $\text{Dist}(\cdot)$ measures the difference between the word embeddings.

This assumption is particularly justified in high-quality pre-trained models, such as CLIP [40], where semantic neighbors in the text space are often mapped to proximity in vector space, supporting a smooth and compact representation of conceptual similarity. However, we admit that while distributional semantics generally ensures semantically similar words are embedded closely, this property varies with different architectures and tokenization strategies. With the assumption above, using $\ell_2$ to stand for $\text{Dist}(\cdot)$, we formalize the target concept as a $\ell_2$-bounded region in the text embedding space. Specifically, for a target concept with text embedding $c_t$, we define its semantic neighborhood as an $\ell_2$ ball:

$$\mathcal{B}(c_t, \gamma) = \{c_t' \in \mathbb{R}^d \mid \|c_t' - c_t\|_2 \leq \gamma\}, \quad (6)$$

where $\gamma$ bounds the maximum distance between $c_t'$ and any synonym embedding. That is, $\mathcal{L}_{align}$ can be further expressed as:

$$\mathbb{E}_{\hat{c}_t \in \mathcal{B}(c_t, \gamma)} \|\hat{c}_t \times (W + \alpha \Delta \hat{W}) - c \times (W + \alpha \Delta W)\|_2^2 \quad (7)$$

We notice that the minimization of this expectation can be transformed into a bi-level optimization problem by maximizing the lower bound of the expectation:

$$\min_{\Delta \hat{W}} \max_{\hat{c}_t \in \mathcal{B}(c_t, \gamma)} \|\hat{c}_t \times (W + \alpha \Delta \hat{W}) - c \times (W + \alpha \Delta W)\|_2^2. \quad (8)$$

However, directly solving this problem often involves large computations since solving the inner maximization is not trivial [34]. Drawing inspiration from recent works [5, 48] illustrating that the adversarial optimization in sample space can also be achieved with adversarially robust parameters within the restricted $\ell_2$ bound:

$$\min_{\Delta \hat{W}} \max_{\|\delta_w\| \leq \tau} \mathbb{E}\|c_t \times (W + \alpha \Delta \hat{W} + \delta_w) - c \times (W + \alpha \Delta W)\|_2^2 \quad (9)$$

Here, $\delta_w$ represents adversarial perturbations to the LoRA parameters, constrained by $\tau$ to balance robustness and stability. Equation 9 shifts the optimization focus from the input embedding space to the parameter space, enabling efficient gradient-based updates while implicitly covering semantic variations within $\mathcal{B}(c_t, \gamma)$.

*4.2.2 Low-Rank Decomposition via Singular Value Decomposition (SVD).* Meanwhile, we should emphasize that directly optimizing $A, B$ to obtain $\Delta \hat{W}$ is not trivial since both matrices are projected into low-dimension space, leading to significantly unstable optimization. Instead, we first initialize $\Delta \hat{W} = BA$ and directly update $\Delta \hat{W}$. Subsequently, $\Delta \hat{W}$ is decomposed via SVD. This step is crucial for reducing the optimization complexity while maintaining the structural integrity of the model's learned representations. Given a weight matrix $\Delta \hat{W}$, we decompose it into two lower-rank matrices $\hat{A}$ and $\hat{B}$, which serve as a compact approximation, ensuring effective pruning alignment with minimal performance loss. Formally, we first perform SVD on $\Delta \hat{W}$:

$$\Delta \hat{W} = U S V^T, \quad (10)$$

where $U$ and $V$ are the left and right singular vector matrices, and $S$ is the diagonal matrix containing singular values. We retain only the top $r$ singular values and their corresponding singular vectors,

---

**Algorithm 1** Robust Editing with Semantic Augmentation

**Require:** Clean T2I diffusion model with text encoder $\mathcal{T}$ and UNet $\epsilon_\theta$; the target concept $C_t$ to be erased; the target LoRA model with parameter $\Delta W$; total editing step $T$.
**Ensure:** Aligned LoRA model.
 1: Initialize $\Delta \hat{W} \leftarrow \Delta W$
 2: $\{c_t^1, c_t^2, c_t^3, \ldots, c_t^K\}, \{c^1, c^2, c^3, \ldots, c^K\} \leftarrow \mathcal{T}(\text{LLM}(C_t))$
 3: **for** $t = 0, 1, \cdots, T - 1$ **do**
 4:     Compute the perturbed $\delta_w$ with $\tau \cdot \nabla_W \mathcal{L}_{align}$
 5:     Solving the inner maximization $\hat{W} \leftarrow W + \delta_w$
 6:     // Compute the loss $\mathcal{L}_{all}$ based on $\hat{W}$.
 7:     $\mathcal{L}_{all} \leftarrow \mathcal{L}_{align} + \eta \cdot \mathcal{L}_{pre}$
 8:     Update $\Delta \hat{W}$ via $\nabla_{\Delta \hat{W}} \mathcal{L}_{all}$ using Adam.
 9: **end for**
10: Approximate $\Delta \hat{W}$ using SVD in Equation 10.
11: **return** $\Delta \hat{W}$

---

indicated as $U_r, S_r, V_r^T$, where $r$ is the rank of LoRA. To construct the final low-rank approximation, we compute the square root of the singular values:

$$\Sigma_r^{\frac{1}{2}} = \text{diag}(\sqrt{S_r}) \quad (11)$$

The resulting decomposition is then defined as:

$$\hat{B} = U_r \Sigma_r^{\frac{1}{2}}, \quad \hat{A} = \Sigma_r^{\frac{1}{2}} V_r^T \quad (12)$$

Thus, the original weight matrix $\Delta \hat{W}$ is approximated as:

$$\Delta \hat{W} \approx \hat{B} \hat{A} \quad (13)$$

Although we only use the approximated version of $\Delta W$, we found that the loss has little influence on the performance of LoRAShield.

### 4.3 Semantic Augmentation

While adversarial optimization allows the erasure of the embedding space around the given concept, it remains confined to local neighborhoods defined by the $\ell_2$ ball $\mathcal{B}(c_t, \gamma)$ and fails to account for global semantic relationships that attackers may exploit through paraphrasing or cross-concept blending. To address this limitation, we propose a semantic augmentation strategy that expands the defended semantic space by enriching the conceptual representation of the target to be erased. Our method operates without requiring access to the original training data, instead leveraging LLMs to generate semantically relevant synonyms and antonyms for the target concept. Formally, given a target concept $c_t$, we query an LLM to produce: (1) Synonyms: Terms closely aligned with $C_t$ (e.g., "modern" → "contemporary") to capture linguistic variations; (2) Antonyms: Terms contrasting with $C_t$ (e.g., "modern" → "archaic") to define conceptual boundaries. For concepts with no viable antonyms, we default to a neutral reference embedding (embedding of ""). This ensures that the erased concept is suppressed without requiring a contrastive counterpart, preserving alignment stability. These augmented terms are then encoded into text embeddings synonyms $\{c_t^1, c_t^2, c_t^3, \ldots, c_t^K\}$ and antonyms $\{c^1, c^2, c^3, \ldots, c^K\}$ ($K = 5$ in this paper) to replace the $c_t$ and $c$ in Equation 9. By integrating these embeddings into the optimization, we enforce robustness

against a broader range of adversarial prompts. The overall procedure of our LoRAShield incorporating "Semantic Augmentation" is illustrated in Algorithm 1.

## 5 EVALUATION

### 5.1 Experimental Setup

*5.1.1 Base Models.* We select SD v1.5 [12] DreamShaper [10] and Realistic Vision [11] as the base models in the experiments, since SD v1.5 is the most popular base model on Civitai and there are many remarkable variants, among which DreamShaper and Realistic Vision, with 1.2 and 1.6 million downloads on Civitai alone [8].

*5.1.2 Datasets and Personalization Methods.* We select four datasets for the evaluation in this paper, including two style-based datasets: (1)"3DM" [31] is a dataset that has 3D rendering style;(2)"pixelart" [46] is a dataset that has pixel-like style; and two portrait-based datasets (3)"Cat" [7] that describe a cat with all kinds of dressings; and (4)"trump" [14] that has the images of Donald Trump generated by LoRA. For LoRA fine-tuning, we set the rank at 128, targeting the attention and projection layer of the UNet with 500 steps (lr=2e-5) following the code given by Hugging Face [25].

*5.1.3 Evaluation Metrics.* Following the previous works [20, 54, 55], the effectiveness of LoRAShield is assessed with the following metrics. (1) Concept Removal Score (CRS) measures the similarity between images generated with prompts with the given word and the target concept using the CLIP score [40], and a lower CRS denotes better erasing performance. Similarly, (2) In-domain Retain Score (IRS) measures the similarity between images generated with prompts of the target concept but different from the given word and the target concept using the CLIP score. (3) Benign Preservation Score (BPS) measures the similarity between the images and the prompts without the undesired concept to evaluate the faithfulness of the image and corresponding prompt using the CLIP score. (4) FID [22] compares the quality and fidelity between the images generated with benign LoRA and images generated with edited LoRA. A lower score means that the distribution of images is more similar. (5)LPIPS compares the semantic similarity between the image generated on benign and the edited LoRA with benign prompts. In particular, for nudity measurement on "3DM" and "pixelart" datasets, following the previous work [27], we use NudeNet [37] to further validate the effectiveness of LoRAShield. The efficiency is measured by running the (1) Time of the erasing process and the (2) Memory costs to conduct the erasing. The robustness of LoRAShield against out-of-domain prompts and attacks is presented in Sec 5.3.

*5.1.4 Implementation Details.* For target concept, following previous settings [19, 27, 56], we select "nude" and as the concept to be erased for style-based datasets (i.e., "3DM" and "pixelart"), and "bloody" for portrait-based datasets (i.e., "cat" and "trump"). For both benign and edited LoRA, we use LLM (Qwen-Plus [50]) to generate 50 prompts with the corresponding trigger word and generate 10 images for each prompt with 10 different random seeds with 500 images for each model and dataset in total. All of the above metrics are evaluated with these generated images. The editing step is set at 10 and τ=1e-5. We set the merging ratio α of LoRA at 1 when merging with the base model in default.

| Prompt | Cat,. . . | 3DM,. . . | 3DM, nude. . . | Cat, bloody. . . |
|---|---|---|---|---|



**Figure 3: Images synthesized with benign (1st row) and edited (2nd row) LoRA when the given prompts without (1st two column) and with (last two column) target concept.**
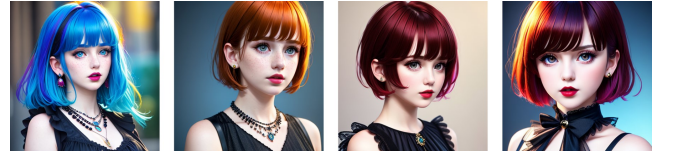


**Figure 4: Cumulative concept erasure: From left to right, images generated by (1) benign LoRA, (2) with concept "hair color" erased, (3) with concepts "hair color"+"necklace" erased, and (4) with concepts "hair color"+"necklace"+"head gesture" erased, illustrating progressive suppression of target semantics while preserving non-target features.**

### 5.2 Performance against Misuse

The performance of LoRAShield is presented in Tab. 2. Note that for a fair comparison, we also report the performance of benign LoRA. To measure the influence of LoRAShield on FID and LPIPS, we regenerate 500 images with different random seeds following the setting before as a comparison for benign LoRAs.

*5.2.1 Erasing Effectiveness.* Tab. 2 illustrates that the images generated by edited LoRA exhibit much lower CLIP scores compared to the benign one, indicating that LoRAShield can erase the target concept effectively. Additionally, when facing in-domain prompts, the edited LoRA can still maintain a low IRS value. Intuitively, we visualize the generated images (on DreamShaper) in Fig. 3, where the images generated on benign LoRA can generate images that significantly violate the protocol of model-sharing platforms. In contrast, even given the explicit prompts, the images generated on edited LoRA present no target concept while aligning with the other semantic prompts. Furthermore, we evaluate the performance of LoRAShield to erase nudity with NudeNet [37] in Tab. 3, in which the images generated on edited LoRA exhibit significantly decrease of nudity and our in-person check further found that many images generated on edited LoRA that were noted as nude by NudeNet, have been misclassified actually.

*5.2.2 Functionality-preserving.* Other than effectiveness, preserving the functionality of LoRA itself is also necessary. As shown in Tab. 2, the BPS values of edited LoRA have a slight decrease

**Table 2: Performance of LoRAShield on three base models and four datasets. Each metric contains two values for edited (1st) and benign (2nd) LoRA, respectively. The ↑ (↓) indicates that a higher (lower) value for the metric signifies superior performance.**

| Base Model | Datasets | Effectiveness | | Functionality-preserving | | | Efficiency | |
|---|---|---|---|---|---|---|---|---|
| | | CRS↓ | IRS↓ | BPS↑ | FID↓ | LPIPS↓ | Time(s) | Memory(GB) |
| **SD1.5** | pixelart | 13.98/17.72 | 11.44/13.22 | 32.06/32.34 | 67.97/94.92 | 0.37/0.67 | 12.97 | 0.2246 |
| | 3DM | 15.24/19.71 | 14.53/17.58 | 30.99/30.56 | 56.44/63.38 | 0.39/0.60 | 13.18 | 0.2363 |
| | trump | 12.70/14.08 | 12.08/13.63 | 25.56/25.27 | 59.60/77.70 | 0.29/0.67 | 13.56 | 0.2246 |
| | Cat | 12.57/14.13 | 12.72/13.42 | 29.56/29.71 | 77.58/69.71 | 0.38/0.62 | 12.2 | 0.2363 |
| **DreamShaper** | pixelart | 12.93/19.30 | 10.68/14.61 | 34.94/35.25 | 70.59/94.34 | 0.41/0.67 | 12.98 | 0.2246 |
| | 3DM | 15.53/20.43 | 14.73/16.90 | 32.47/32.70 | 40.36/54.88 | 0.30/0.54 | 13.74 | 0.2246 |
| | trump | 11.34/14.34 | 11.67/14.41 | 26.29/26.35 | 66.21/68.29 | 0.40/0.62 | 17.94 | 0.2178 |
| | Cat | 12.83/14.62 | 12.50/13.56 | 30.90/31.01 | 68.02/65.21 | 0.37/0.56 | 16.32 | 0.2363 |
| **Realistic Vision** | pixelart | 13.20/20.41 | 10.88/14.56 | 35.34/35.01 | 65.27/108.55 | 0.37/0.67 | 13.52 | 0.2246 |
| | 3DM | 16.55/20.38 | 14.73/17.35 | 31.89/32.44 | 46.50/58.88 | 0.32/0.55 | 11.68 | 0.2363 |
| | trump | 12.40/14.81 | 12.46/13.56 | 26.22/26.78 | 62.00/69.69 | 0.35/0.61 | 17.05 | 0.2168 |
| | Cat | 12.75/14.40 | 12.05/13.30 | 31.21/30.84 | 66.62/65.69 | 0.35/0.57 | 13.75 | 0.2207 |

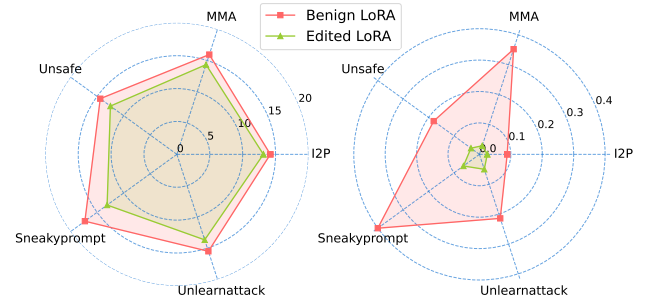**Table 3: Nudity of images generated with edited and benign LoRA on three base models and two datasets.**

| Base Model | Dataset | Edited | Benign |
|---|---|---|---|
| SD1.5 | 3DM | 0.09 | 0.67 |
| SD1.5 | pixelart | 0.06 | 0.28 |
| DreamShaper | 3DM | 0.02 | 0.63 |
| DreamShaper | pixelart | 0.09 | 0.55 |
| Realistic | 3DM | 0.13 | 0.74 |
| Realistic | pixelart | 0.09 | 0.59 |

**Table 4: Performance against multiple LoRA merge.**

| LoRA | Benign | | Edited | |
|---|---|---|---|---|
| | CRS | Nudity | CRS | Nudity |
| Base | 20.43 | 0.63 | 15.53 | 0.02 |
| +Tools [13] | 20.60 | 0.70 | 15.36 | 0.08 |
| +Background [36] | 16.39 | 0.22 | 14.11 | 0.02 |
| +Clothes [28] | 15.11 | 0.10 | 13.77 | 0.06 |
| +Pose [1] | 14.99 | 0.06 | 13.64 | 0.00 |
| +Celebrity [3] | 13.81 | 0.06 | 12.72 | 0.00 |

in some cases, but most keep close to that of benign LORA. This can be attributed to the interrelationship of concepts, and erasing the target concept may affect the representation of another one. Surprisingly, the images generated on edited LoRA have lower FID and LPIPS values than images generated on benign ones (with different random seeds), which also validates that the negative impact of LoRAShield is acceptable, since its influence is much less than the change of random seeds. From the visual comparison of the images generated on the edited and benign LoRAs in Fig. 3, we can conclude that LoRAShield can effectively preserve the core functionality of LoRAs, aligning with our design goal of balancing robustness with minimal disruption to legitimate generation.

*5.2.3 Erasing Efficiency.* LoRAShield achieves platform-scale efficiency, which is critical for deployment in model-sharing ecosystems. As shown in Tab. 2, editing a LoRA takes 14 seconds (vs. often several minutes for retraining from scratch) with 0.23GB GPU memory. This efficiency stems from our sequential optimization of attention matrices, which avoids full model retraining and leverages LoRA's low-rank structure. Even for platforms hosting 100k+ LoRAs, LoRAShield enables fast processing on a single A100 GPU, ensuring real-time user experiences even during peak uploads.



**Figure 5: CRS (left) and Nudity (right) under input attacks.**

*5.2.4 Multiple Concept Erasure.* To evaluate LoRAShield's ability to handle cumulative semantic suppression, we conducted a dynamic visualization experiment erasing up to three interdependent concepts in Fig. 4, since quantitative metrics (e.g., CLIP scores) become less reliable for multi-concept scenarios.

## 5.3 Effectiveness against Adversarial Attack

As mentioned before in Section 3.3.3, since malicious users can download the edited LoRA from the platforms, it is critical to validate whether LoRAShield can resist potential attacks. In this paper, we consider four input-based attacks [39, 42, 51, 52] that construct adversarial prompts. Another common scenario is that malicious users may merge multiple LoRAs for misuse, which can disrupt the protection of LoRAShield. For simplicity, we consider "nude" as the target concept on "3DM" datasets and use "DreamShaper" as base model. More details can be found in the appendix.

*5.3.1 Robustness against Multi-LoRA Merge.* A common but critical threat arises when attackers merge multiple edited LoRAs. For instance, combining a "beach scenery" LoRA (with "nude" erased) and "human portrait" LoRA (without "nude" erased) might reintroduce NSFW content. We use weight averaging to linearly combine multiple LoRA matrices $\Delta W_{merged} = \sum \alpha_i \Delta W_i$. Five popular LoRAs of different types from Civitai: tools [13] (for adjusting the clothes of people), background [36], clothes [28], pose [1] and celebrity [3] are selected for evaluation with base LoRA trained on "3DM" dataset and DreamShaper with $\alpha_i = 1$. The result in Tab. 4 showcases that even under the perturbation in parameter space, the edited LoRA retains stable performance compared with the base one.

*5.3.2 Robustness against Input-based Attack.* Malicious users may iteratively refine the prompts to bypass concept erasure, escalating from simple keyword swaps to sophisticated semantic manipulations. To simulate this arms race, we evaluate LoRAShield against four SOTA input-based attacks [39, 42, 51, 52] that reflect realistic evasion tactics. The result in Tab. 5 demonstrates that even under input-based attacks, LoRAShield exhibits remarkable performance. Specifically, the nudity score of images generated by edited LoRA is much lower than that of benign one.

## 5.4 Real-world Case Study

To validate the practical efficacy of LoRAShield, we conducted end-to-end experiments on Civitai. Our threat scenario mirrors real-world misuse: benign LoRAs [43] uploaded by legitimate users are exploited by attackers to generate policy-violating content, as shown in the first row of Fig. 6. We can notice that the benign LoRA flagged for generating celebrity portraits but could produce NSFW content when given specified prompt. To mitigate this, we download it and attempt to generate NSFW images with the same prompt on edited LoRA. The second row of Fig. 6 further highlights that LoRAShield empowers platforms to proactively neutralize misuse risks without sacrificing creative utility or scalability.

## 5.5 Ablation Study

To dissect LoRAShield's design choices, we performed ablation experiments on the following critical factors: robust optimization, semantic augmentation, and merging ratio. All experiments in this subsection were conducted on the "3DM" dataset and DreamShaper.

*5.5.1 Impact of Robust Optimization and Semantic Augmentation.* Since we have claimed that the proposed robust optimization and semantic augmentation enhance the generalization of concept erasure, we validate this with an ablation study. The result in Tab. 5
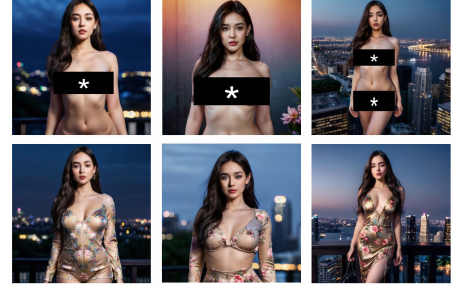


**Figure 6: Generated images using LoRA [43] downloaded on Civitai. Note that these images will never be posted online.**

**Table 5: Ablation study on robust optimization (RO) and semantic augmentation (SA). ✔ and ✘ denote whether the corresponding techniques are incorporated in LoRAShield.**

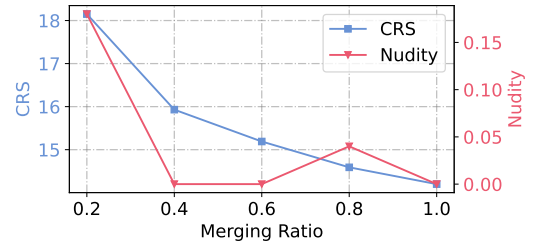| RO | SA | CRS | Nudity |
|----|----|-----|--------|
| ✘ | ✘ | 14.60 | 0.02 |
| ✘ | ✔ | 14.55 | 0.00 |
| ✔ | ✘ | 14.53 | 0.00 |
| ✔ | ✔ | 14.20 | 0.00 |



**Figure 7: Impact of Merging Ratio.**

illustrates that both enhancements can decrease the nudity score of the generated images while incorporating together further improves CRS, validating our design of LoRAShield.

*5.5.2 Impact of Merging Ratio.* The merge $\alpha$ ratio governs the strength of LoRA weights when merging with the base model, and a lower strength may weaken the protection of LoRAShield as well as the benign concept. Here, we set $\alpha$ at {0.2, 0.4, 0.6, 0.8, 1.0} to quantify trade-offs between faithfulness and protection and the result in Fig. 7 demonstrates that even the nudity score fluctuates slightly but remains at a low level when $\alpha \geq 0.4$. Importantly, we should also emphasize that a low merging ratio will definitely degrade the performance of benign tasks. In real-world misuse scenarios, attackers are incentivized to maintain high $\alpha$ values to preserve the LoRA's utility for legitimate tasks while attempting to bypass safeguards. This constraint ensures that the protection remains effective when adversaries prioritize functional performance.

# 6 CONCLUSION

In this paper, we identify critical vulnerabilities in model-sharing ecosystems. While DMs have been extensively studied for concept erasure, LoRA's misuse risk has been largely overlooked, leaving platforms like Civitai exposed to misuse that harms creators' rights and platform integrity. This work bridges this gap by proposing LoRAShield, the first data-free and editing framework to secure LoRAs against adversarial exploitation. By editing LoRA's weight subspace, our platform-driven approach achieves effectiveness and efficiency in suppressing harmful content while preserving functionality for legitimate uses, enabling scalable deployment across millions of shared LoRAs.

# REFERENCES

[1] Akiseki. 2024. Shirt tug pose lora. https://civitai.com/models/7706/shirt-tug-pose-lora
[2] Lucas Bourtoule, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. 2021. Machine unlearning. In *2021 IEEE symposium on security and privacy (SP)*. IEEE, 141–159.
[3] Bravenunu. 2024. MiMi. https://civitai.com/models/14816/mimi
[4] Heather Chen and Kathleen Magramo. 2024. Finance worker pays out $25 million after video call with deepfake 'chief financial officer'. https://edition.cnn.com/2024/02/04/asia/deepfake-cfo-scam-hong-kong-intl-hnk
[5] Jiahao Chen, Zhou Feng, Rui Zeng, Yuwen Pu, Chunyi Zhou, Yi Jiang, Yuyou Gan, Jinbao Li, and Shouling Ji. 2024. Enhancing Adversarial Transferability with Adversarial Weight Tuning. *arXiv preprint arXiv:2408.09469* (2024).
[6] Zhi-Yi Chin, Chieh-Ming Jiang, Ching-Chun Huang, Pin-Yu Chen, and Wei-Chen Chiu. 2023. Prompting4debugging: Red-teaming text-to-image diffusion models by finding problematic prompts. *arXiv preprint arXiv:2309.06135* (2023).
[7] chosen. 2025. Cute cat midjourney style cat lora. https://civitai.com/models/99579/cute-cat-midjourney-style-cat-lora?modelVersionId=106565
[8] Inc. Civit AI. 2025. Civitai Models. https://civitai.com/models
[9] Inc. Civit AI. 2025. Civitai Safety Center. https://civitai.com/safety
[10] Inc. Civit AI. 2025. DreamShaper. https://civitai.com/models/4384/dreamshaper
[11] Inc. Civit AI. 2025. Realistic Vision. https://civitai.com/models/4201/realistic-vision-v60-b1
[12] Inc. Civit AI. 2025. Stable-Diffusion-1.5. https://huggingface.co/Jiali/stable-diffusion-1.5
[13] Leosams clothing adjuster lora. 2025. LEOSAM. https://civitai.com/models/88132/leosams-clothing-adjuster-lora
[14] coplasx. 2025. donald-trump. https://civitai.com/models/17171/donald-trump
[15] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems* 34 (2021), 8780–8794.
[16] Martin Nicolas Everaert, Athanasios Fitsios, Marco Bocchio, Sami Arpa, Sabine Süsstrunk, and Radhakrishna Achanta. 2024. Exploiting the signal-leak bias in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 4025–4034.
[17] Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. 2023. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. *arXiv preprint arXiv:2310.12508* (2023).
[18] fxc2017gd635. 2024. Pure Innocent Girl. https://civitai.com/models/35685/pure-innocent-girl?modelVersionId=78583
[19] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. 2023. Erasing concepts from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2426–2436.
[20] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. 2024. Unified concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 5111–5120.
[21] Jason Griffin. 2024. Civitai: The Challenges of Fairly Monetising a Generative AI Community. https://www.linkedin.com/pulse/civitai-challenges-fairly-monetising-generative-ai-jason-griffin-usuxe
[22] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* 30 (2017).
[23] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR* 1, 2 (2022), 3.

[24] HuggingFace. 2025. The AI community building the future. https://huggingface.co/
[25] HuggingFace. 2025. Train text to image lora. https://github.com/huggingface/diffusers/blob/main/examples/text_to_image/train_text_to_image_lora.py
[26] Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu. 2023. Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 22691–22702.
[27] Xinfeng Li, Yuchen Yang, Jiangyi Deng, Chen Yan, Yanjiao Chen, Xiaoyu Ji, and Wenyuan Xu. 2024. Safegen: Mitigating sexually explicit content generation in text-to-image models. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*. 4807–4821.
[28] liaoliaojun. 2024. hanfu. https://civitai.com/models/15365/hanfu
[29] liblib. 2025. LiblibAI. https://www.liblib.art/
[30] LiblibAI. 2024. User Agreement. https://www.liblib.art/us/user
[31] LONGD. 2025. 3d-rendering-style. https://civitai.com/models/73756/3d-rendering-style
[32] Sasha Luccioni, Christopher Akiki, Margaret Mitchell, and Yacine Jernite. 2023. Stable bias: Evaluating societal representations in diffusion models. *Advances in Neural Information Processing Systems* 36 (2023), 56338–56351.
[33] Mengyao Lyu, Yuhong Yang, Haiwen Hong, Hui Chen, Xuan Jin, Yuan He, Hui Xue, Jungong Han, and Guiguang Ding. 2024. One-dimensional adapter to rule them all: Concepts diffusion models and erasing applications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7559–7568.
[34] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083* (2017).
[35] malcolmrey. 2024. Joe Biden. https://civitai.com/models/533408/joe-biden
[36] Mccc. 2023. Ghibli background. https://civitai.com/models/54233/ghiblibackground
[37] notAI tech. 2024. Lightweight nudity detection. https://github.com/notAI-tech/NudeNet
[38] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2023. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952* (2023).
[39] Yiting Qu, Xinyue Shen, Xinlei He, Michael Backes, Savvas Zannettou, and Yang Zhang. 2023. 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*. 3403–3417.
[40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PmLR, 8748–8763.
[41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
[42] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. 2023. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 22522–22531.
[43] seanwang1221. 2024. Li Yitong CN actress SD15 & FLUX. https://civitai.com/models/44324/li-yitong-cn-actress-sd15-and-flux?modelVersionId=82580
[44] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2023. Understanding and mitigating copying in diffusion models. *Advances in Neural Information Processing Systems* 36 (2023), 47783–47803.
[45] 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. [n. d.]. Ring-A-Bell! How Reliable are Concept Removal Methods For Diffusion Models?. In *The Twelfth International Conference on Learning Representations*.
[46] VatsaDev. 2025. Train text to image lora. https://huggingface.co/datasets/VatsaDev/pixel-art
[47] Yiluo Wei, Yiming Zhu, Pan Hui, and Gareth Tyson. 2024. Exploring the Use of Abusive Generative AI Models on Civitai. In *Proceedings of the 32nd ACM International Conference on Multimedia*. 6949–6958.
[48] Dongxian Wu, Shu-Tao Xia, and Yisen Wang. 2020. Adversarial weight perturbation helps robust generalization. *Advances in neural information processing systems* 33 (2020), 2958–2969.
[49] Jing Wu, Trung Le, Munawar Hayat, and Mehrtash Harandi. 2024. Erasediff: Erasing data influence in diffusion models. *arXiv preprint arXiv:2401.05779* (2024).
[50] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115* (2024).
[51] Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Tsung-Yi Ho, Nan Xu, and Qiang Xu. 2024. Mma-diffusion: Multimodal attack on diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7737–7746.

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

[53] Gong Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. 2024. 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*. 1755–1764.

[54] Yimeng Zhang, Xin Chen, Jinghan Jia, Yihua Zhang, Chongyu Fan, Jiancheng Liu, Mingyi Hong, Ke Ding, and Sijia Liu. 2024. Defensive unlearning with adversarial training for robust concept erasure in diffusion models. *Advances in Neural Information Processing Systems* 37 (2024), 36748–36776.

[55] Yihua Zhang, Chongyu Fan, Yimeng Zhang, Yuguang Yao, Jinghan Jia, Jiancheng Liu, Gaoyuan Zhang, Gaowen Liu, Ramana Rao Kompella, Xiaoming Liu, et al.

[n. d.]. UnlearnCanvas: Stylized Image Dataset for Enhanced Machine Unlearning Evaluation in Diffusion Models. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

[56] Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and Sijia Liu. 2024. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. In *European Conference on Computer Vision*. Springer, 385–403.

[57] Yihua Zhang, Yimeng Zhang, Yuguang Yao, Jinghan Jia, Jiancheng Liu, Xiaoming Liu, and Sijia Liu. 2024. Unlearncanvas: A stylized image dataset to benchmark machine unlearning for diffusion models. *arXiv e-prints* (2024), arXiv–2402.

[58] Xin Zhao, Xiaojun Chen, Yuexin Xuan, Zhendong Zhao, Xiaojun Jia, Xinfeng Li, and Xiaofeng Wang. 2024. Buster: Implanting Semantic Backdoor into Text Encoder to Mitigate NSFW Content Generation. *arXiv preprint arXiv:2412.07249* (2024).