## THE DUAL POWER OF INTERPRETABLE TOKEN EMBED-DINGS: JAILBREAKING ATTACKS AND DEFENSES FOR DIFFUSION MODEL UNLEARNING

## **Anonymous authors**

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

035

036

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

#### **ABSTRACT**

Diffusion models excel at generating high-quality images but can memorize and reproduce harmful concepts when prompted. Although fine-tuning methods have been proposed to unlearn a target concept, they struggle to fully erase the concept while maintaining generation quality on other concepts, leaving models vulnerable to jailbreak attacks. Existing jailbreak methods demonstrate this vulnerability but offer limited insight into how unlearned models retain harmful concepts, limiting progress on effective defenses. In this work, we take one step forward by exploring a linearly interpretable structure. We introduce SubAttack, a novel jailbreaking attack that learns an orthogonal set of attack token embeddings, each being a linear combination of human-interpretable textual elements, revealing that unlearned models still retain the target concept through related textual components. Furthermore, our attack is also more powerful and transferable across text prompts, initial noises, and unlearned models than prior attacks. Leveraging these insights, we further propose SubDefense, a lightweight plug-and-play defense mechanism that suppresses the residual concept in unlearned models. SubDefense provides stronger robustness than existing defenses while better preserving safe generation quality. Extensive experiments across multiple unlearning methods, concepts, and attack types demonstrate that our approach advances both understanding and mitigation of vulnerabilities in diffusion unlearning.

## 1 Introduction

Diffusion models (DMs) have recently emerged as a powerful class of generative models, capable of producing diverse and high-quality content such as images (Ho et al., 2020), videos (Khachatryan et al., 2023), and protein structures (Watson et al., 2023). Notably, Text-to-Image (T2I) diffusion models (Rombach et al., 2022; Ramesh et al., 2022a; Saharia et al., 2022; Zhang et al., 2024e;b) have gained significant popularity for their ability to generate high-fidelity images from user-provided text prompts. However, the remarkable generative capabilities of these models also raise significant concerns regarding their safe deployment. For example, users can exploit carefully crafted text prompts to induce these models by generating unethical or harmful content, such as nude or violent images, or copyrighted material (Schramowski et al., 2023).

To address such safety concerns without refiltering the huge dataset and retraining the full model, *Machine Unlearning* (MU) methods have recently been developed for "erasing" a harmful concept directly from the pretrained models. For instance, a wide range of methods (Gandikota et al., 2023; 2024; Zhang et al., 2023; Lyu et al., 2024) seek to unlearn harmful content in pretrained DMs by finetuning the model weights (Nguyen et al., 2024). Yet, the key challenge of preserving the generation quality of safe content limits unlearned DMs from removing even a *single concept* completely. This limitation becomes evident under *jailbreaking attacks* (Zhang et al., 2024d; Pham et al., 2024; Chin et al., 2024b; Tsai et al., 2024; Zhuang et al., 2023), which have enforced unlearned DMs to regenerate harmful content. For instance, UnlearnDiff (Zhang et al., 2024d) crafts adversarial discrete text prompts, and CCE (Pham et al., 2024) leverages textual inversion (Gal et al., 2023) to execute jailbreaking attacks in embedding space. Amid the rising popularity of open-source models, and given the risks of insider threats and model leakage, many studies adopt a white-box setting (i.e., full access to model weights) for safety evaluation. These works reveal that unlearned DMs remain vulnerable, highlighting the urgent need to *defend* them by strengthening robustness against attacks.

It is not surprising that optimization-based, non-interpretable, and worst-case prompt perturbations can jailbreak unlearned DMs. However, despite leveraging white-box access, such approaches provide limited *interpretability*, i.e., a human-understandable explanation of how a model's internal state drives the prediction of its behavior under intervention. Therefore, they offer little insight into how harmful concepts persist within the model, and these attacks fail to offer potential insights for defense strategies. Furthermore, the defense of unlearned DMs remains largely underexplored. For instance, the RECE defense framework (Gong et al., 2024) focuses on improving a specific unlearned model (UCE (Gandikota et al., 2024)) against particularly adversarial attacks (i.e., UnlearnDiff). Extending defenses to a broader range of unlearned models and attack types remains a challenging problem. These gaps motivate our central **question**: *Can we design more human-interpretable jailbreaking attacks that also provide actionable insights for building defenses in unlearned DMs?* 

Our work tackles this fundamental question by exploring underlying linear structures, taking advantage of the white-box setting. We introduce an effective, human-interpretable *subspace attack method (SubAttack)*, which further inspires a *subspace defense strategy (SubDefense)* broadly applicable to various unlearned models and attacks. The core idea is to learn an orthogonal set of attack token embeddings within the unlearned model for the harmful concept. Inspired by prior works (Chefer et al., 2024; Park et al., 2023a; Cunningham et al., 2023), we optimize each attack embedding as a nonnegative linear combination of embeddings of existing concepts, and interpret the concept through the linear decomposition. Leveraging our approach, we show that unlearned DMs associate the harmful concept with mixtures of other hidden concepts, thus retaining unintended harmful regeneration capabilities. These insights motivate our defense mechanism, which further mitigates the harmful concept from unlearned DMs by removing the learned attack token embeddings through orthogonal subspace projection.

Compared to prior methods, our SubAttack demonstrates strong empirical performance of efficiency and effectiveness, showing stronger transferability across text prompts, initial noises, and unlearned models. Our defense strategy can be seamlessly integrated into various unlearned models, improving robustness against different jailbreaking attacks while preserving higher generation quality than the baseline defense method (Gong et al., 2024). A comprehensive discussion of related works is in **App. A**. In summary, this work makes the following **contributions**:

- **Interpretable attack via linear structure.** We propose *SubAttack*, which learns an orthogonal set of token embeddings under a linear structure. These embeddings can be interpreted in a bag-of-words fashion, revealing how the residual concept is still retained in unlearned DMs.
- Effective and transferable attack. *SubAttack* achieves higher ASR than existing baselines across diverse concepts and unlearned models, while also transferring reliably across prompts, initial noise, and models, exposing critical vulnerabilities in current unlearning methods.
- Subspace defense inspired by the attack. Leveraging this linear structure, we propose *SubDefense*, which projects out attack token directions to eliminate residual concepts. Our SubDefense offers versatile, reliable protection while preserving generation quality.

## 2 Preliminaries and Problem Statement

## 2.1 Preliminaries

Overview of Latent Diffusion Models (LDMs). T2I diffusion models have recently gained popularity for their ability to generate desired images from user-provided text prompts. Among these various T2I models, LDMs (Rombach et al., 2022) is the most widely deployed DM, and has therefore become the primary focus of current machine unlearning methods. As shown in Fig. 1, for a given text prompt p, LDM first encodes p using a pretrained CLIP text encoder (Radford et al., 2021)  $f(\cdot)$  to obtain the text embedding c = f(p). Then, the generation process begins by sampling a random noise  $z_T \sim \mathcal{N}(0,1)$  in the latent space. After that, LDM progressively denoises  $z_T$  conditioned on the context c until the final clean latent  $z_0$  is achieved. Specifically, for each timestep  $t = T, T - 1, \ldots, 1$ , its denoising UNet,  $\epsilon_{\theta}(z_t \mid c)$ , predicts and removes the noise to obtain a cleaner latent representation  $z_{t-1}$ . The clean latent  $z_0$  is then decoded to an image with a pretrained image decoder. To train the denoising UNet  $\epsilon_{\theta}(z_t \mid c)$  in LDM, the denoising error is minimized:

$$\mathcal{L} = \mathbb{E}_{(\boldsymbol{z},\boldsymbol{c}),t,\epsilon \sim \mathcal{N}(0,1)} \left[ \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}} \left( \boldsymbol{z}_t \mid \boldsymbol{c} \right) \right\|_2^2 \right], \tag{1}$$

where z is the clean image latent encoded by a pretrained image encoder and c is the corresponding text embedding. Here,  $z_t = \sqrt{\alpha_t}z + \sqrt{1-\alpha_t}\epsilon$  is the noisy image latent at timestep t, and  $\alpha_t > 0$  is a predefined constant.

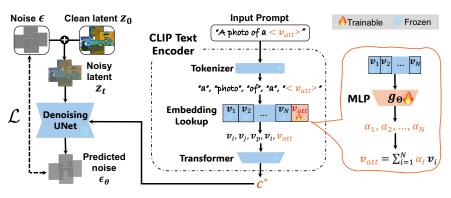


Figure 1: Learning one interpretable attack token embedding. The learning process of one attack token embedding  $v_{\rm att}$  for the concept "Van Gogh" is visualized. Blue parts represent the frozen unlearned LDM, where, for simplicity, we omit the image encoder and decoder. In orange parts, it illustrates the learning mechanism for optimizing an MLP network to produce  $v_{\rm att}$ , which is a linear combination of the existing token embeddings.

**CLIP text encoder and the token embedding space.** To control the generation process, a key component of LDMs is the pretrained CLIP text encoder  $f(\cdot)$ . As illustrated in **Fig. 1**, the CLIP text encoder consists of three main components:

- **Tokenizer:** This module splits the text prompt *p* into a sequence of tokens, which can be words, sub-words, or punctuation marks. Each token is assigned a unique token ID from the CLIP text encoder's predefined vocabulary.
- Token Embeddings: These token IDs (e.g.,  $[i, j, \cdots]$ ) are then mapped to corresponding token embeddings  $v_i \in \mathbb{R}^d$  stored in the token embedding table. This process generates a sequence of token embeddings  $[v_i, v_j, \cdots]$ .
- Transformer Network: This network processes the sequence of token embeddings and encodes them into the final text embedding c that can guide the image generation process in LDMs.

Through optimizing Eq. (1), LDM learns to associate activations in the text encoder with concepts in the generated images. Prior research has explored controlling generated content through manipulating activations in the text encoder. In particular, it has been identified that the token embedding space v plays a vital role in content personalization, where a single text embedding can represent a specific attribute (Gal et al., 2023) and the token embedding space can be utilized for linear decomposition of concepts (Chefer et al., 2024). Leveraging the expressiveness and interpretability of the token embedding space, we propose both jailbreaking attack and defense mechanisms, and discuss the problem setup in the following.

## 2.2 PROBLEM STATEMENT AND SETUP

Jailbreaking attacks are designed to evaluate the robustness of unlearned LDMs. Most existing diffusion unlearning studies focus on removing a single target concept from each model. For example, given a prompt p= "a photo of a [target concept] ...", an unlearned LDM for this concept is expected to have difficulty generating the corresponding images. A jailbreaking attack, given an unlearned LDM as the victim model, aims to manipulate the prompt to make the model regenerate the unwanted concept. There are majorly two kinds of attack setups: (i) Adversarial prompt-based attacks (Zhang et al., 2024d; Chin et al., 2024b; Tsai et al., 2024; Zhuang et al., 2023) optimize an adversarial text prompt  $p_{\text{att}}$  and append it to p. (ii) Embedding-based attacks (Pham et al., 2024) learn an attack token embedding  $v_{\text{att}}$ , register it as a new token  $< v_{\text{att}} >$ , and modify the prompt by replacing the [target concept] with this token. Our attack follows the second setup, but is explicitly designed to be interpretable through linear constraints while achieving stronger attack performance. Moreover, apart from access to the unlearned LDM, existing attacks generally require either the original LDM or images containing the target concept; in our setup, we assume access to the images ( $z_0$  in Fig. 1).

**Defense**, in contrast, seeks to protect an *unlearned LDM* from new jailbreaking attacks. Once a defense strategy is applied, it should prevent the model from regenerating harmful concepts even under *future attacks*, while preserving its ability to generate harmless content. For example, RECE (Gong et al., 2024) further modifies the denoising UNet of the unlearned model UCE (Gandikota et al., 2024) to defend against adversarial attacks (Zhang et al., 2024d). In this work, we propose a

163

164

165

166

167

169

170 171

172

173

174

175

176

177

178

179

181

182

183

185

186

187 188 189

190 191

192

193

195

196

197

199

200

201202

203

204

205206

207

208209

210

211

212

213

214

215

Figure 2: Interpreting the attack token embeddings for concept "nudity", "Van Gogh", and "church". Tokens with the largest  $\alpha_i$  are words associated with the target concept. For example, top tokens for "church" are activities conducted in the church, or names from the Bible.

defense strategy that safeguards the token embedding space and can be seamlessly integrated into existing unlearned LDMs. Our objective is to provide a broadly applicable defense mechanism that enhances robustness across diverse unlearned models when confronted with new attacks.

Interpretability refers to providing a compact, human-understandable description of how a model's internal components drive its behavior (Chefer et al., 2024; Zou et al., 2023; Bereska & Gavves, 2024). Such a description is crucial as it allows testable predictions under controlled interventions. In our setting, interpretability means that attack embeddings can be explained as recognizable words or semantic units rather than opaque vectors. While many existing methods are empirically effective, they lack such interpretability, making it difficult to understand how harmful concepts persist or how to control the robustness of unlearned models. Our work addresses this gap by developing attack and defense methods that are both more interpretable and effective.

**Notations.** Before introducing our method, we define the following projection operators. Specifically, given vector z, for a vector v, let  $\operatorname{Proj}_v(z)$  denote the projection of z onto v. For a matrix V, let  $\operatorname{Proj}_V(z)$  denote the projection of z onto the subspace spanned by the columns of V. Formally, these operators are given by

$$\operatorname{Proj}_{oldsymbol{v}}(oldsymbol{z}) := rac{oldsymbol{v} oldsymbol{v}^ op}{\|oldsymbol{v}\|_2^2} oldsymbol{z}, \ \operatorname{Proj}_{oldsymbol{V}}(oldsymbol{z}) := oldsymbol{V}(oldsymbol{V}^ op oldsymbol{V})^{-1} oldsymbol{V}^ op oldsymbol{z}.$$

## 3 SUBSPACE ATTACKING AND DEFENDING METHODS

This section introduces our subspace attacking and defending methods for LDMs. In Sec. 3.1, we explore the token embedding space to develop an interpretable and effective attack method (SubAttack) by learning a sequence of attack token embeddings orthogonal to each other. SubAttack further inspires us to propose a defense strategy (SubDefense) in Sec. 3.2, by orthogonal subspace projection of learned attack token embeddings, which can effectively defend against various jailbreaking attacks.

#### 3.1 SUBSPACE ATTACKING: SubAttack

Before we introduce our subspace attacking method, let us build some intuition of how to learn a single interpretable attack token embedding  $v_{\text{att}} \in \mathbb{R}^d$ . Based on this, we will then show how to iteratively learn a sequence of orthogonal attack token embeddings through *deflation*, i.e., removing already computed embeddings.

## 3.1.1 SINGLE-TOKEN EMBEDDING ATTACK

We aim to learn a token embedding  $v_{att} \in \mathbb{R}^d$  as a non-negative linear representation of existing token embeddings  $v_i$  in the CLIP vocabulary  $\mathcal{V}$  as follows:

$$v_{\text{att}} = \sum_{i=1}^{N} \alpha_i v_i, \quad \alpha_i = g_{\Theta}(v_i) \ge 0,$$
 (2)

where N is the total size of the original CLIP vocabulary, and  $v_i$ ,  $i=1,2,\ldots,N$ , are original CLIP token embeddings within  $\mathcal{V}$ . Non-negative  $\alpha_i$  are parameterized via a multi-layer perceptron (MLP) network  $g_{\Theta}(\cdot): \mathbb{R}^d \mapsto \mathbb{R}^+$  with ReLU activation. This is inspired by recent work (Chefer et al., 2024) on language models. To learn  $v_{\text{att}}$ , we **optimize the loss**  $\mathcal{L}$  in Eq. (1) with respect to the **parameter**  $\Theta$  **of the MLP**, while freezing all the other components. As illustrated in **Fig. 1**, during training we enforce the training data pairs  $(z, c^*) \sim \mathcal{D}$  to satisfy the following constraints: (i) z is the latent image containing the target harmful concept. (ii)  $c^*$  is the text embedding for the text prompt p, and p contains the new special token  $< v_{\text{att}} >$  whose token embedding is  $v_{\text{att}}$ .

**Remarks.** The non-negative constraint in Eq. (2) is inspired by prior works on linear representation hypothesis and linear feature decomposition (Chefer et al., 2024; Zhou et al., 2018; Cunningham et al., 2023; Park et al., 2023a) that "negative concepts are not as interpretable as positive concepts." In this way, the target concept can be viewed as a combination of top-weighted (i.e., having largest  $\alpha_i$ ) concepts in  $\mathcal{V}$ . **Fig. 2** illustrates the identified sets of human-interpretable concepts for different target concepts (e.g., nudity, Van Gogh, church) in unlearned LDMs. Additionally, we provide analysis on the sparsity of  $\alpha_i$  in **App. F**. Now, we introduce how a set of attack token embeddings are learned.

#### 3.1.2 Subspace Token Embedding Attacks

Compared with learning a single attack token embedding  $v_{\text{att}}$ , it is more powerful to learn a set of diverse attacks  $\{v_{\text{att},k}\}_{k=1}^K \ (m \leq d)$  that can attack the same target concept, as outlined in **Algorithm 1**. We enforce orthogonality on  $\{v_{\text{att},k}\}_{k=1}^K$  to promote diversity and improve attack effectiveness (see ablations in **App. E.1**).

Such a set of orthogonal token embeddings  $\{v_{\text{att},k}\}_{k=1}^K$  is learned through deflation, sharing similar spirits with classical numerical methods such as orthogonal matching pursuit (Tropp & Gilbert, 2007). Specifically, suppose the first attack token embedding  $v_{\text{att,1}}$  is identified following Sec. 3.1.1 by optimizing an MLP  $g_{\Theta_1}$ ,

```
Algorithm 1 Learning Attack Token Embeddings

1: Input: victim model with CLIP token embeddings [v_{1,1},\ldots,v_{N,1}], total iterations K

2: Output: [v_{\mathsf{att},1},\ldots,v_{\mathsf{att},K}]

3: for k=1,\ldots,K do

4: Optimize the MLP g_{\Theta_k}

5: \alpha_{i,k} \leftarrow g_{\Theta_k}(v_{i,k})

6: v_{\mathsf{att},k} \leftarrow \sum_{i=1}^N \alpha_{i,k} v_{i,k}

7: for i=1,\ldots,N-1 do

8: v_{i,k+1} \leftarrow v_{i,k} - \operatorname{Proj}_{v_{\mathsf{att},k}}(v_{i,k})

9: end for
```

we then "eliminate" the target concept  $v_{{\tt att},1}$  from the whole vocabulary  ${\cal V}$  via orthogonal projection:

10: **end for** 

$$\mathbf{v}_{i,2} = \mathbf{v}_{i,1} - \operatorname{Proj}_{\mathbf{v}_{\text{att},1}}(\mathbf{v}_{i,1}), \quad \forall \ i \in [N].$$
 (3)

Here,  $v_{i,1} \equiv v_i \in \mathcal{V}$  are the original embeddings for all  $i \in [N]$ . Eq. (3) makes sure all the updated  $v_{2,i},\ldots,v_{2,N}$  are orthogonal to  $v_{\text{att},1}$ . With the new  $\mathcal{V}_2 = \{v_{2,i}\}_{i=1}^N$ , we can learn a second attack token embedding  $v_{\text{att},2} = \sum_{i=1}^N \alpha_{i,2} v_{i,2}, \ \alpha_{i,2} = g_{\Theta_2}(v_{i,2}) \geq 0$ , then  $v_{\text{att},2}$  is ensured to be orthogonal to  $v_{\text{att},1}$ . Here,  $g_{\Theta_2}$  is another MLP optimized in the same way as  $g_{\Theta_1}$ . As such, we can repeat the procedure for K times to learn and construct a set of orthogonal attack token embeddings  $\{v_{\text{att},k}\}_{k=1}^K$ , and use each of them to attack the same target concept. In practice, during attacking, we choose K=5, which delivers strong attack performance while keeping the method efficient (see ablation studies in App. E.1).

## 3.2 Subspace Defending: SubDefense

Our SubAttack reveals that combinations of related hidden concepts can represent the target concept in an unlearned LDM through a linear composition. This insight motivates us to design a defense strategy within the same linear framework. Our intuition is to remove these identified concept representations from unlearned models through orthogonal projection, thereby making them more robust to various jailbreaking attacks. Concretely, because linearly composed concepts become more difficult to recover, this is achieved by projecting onto the null space of the learned subspace attacks.

More specifically, suppose we have learned a set of attack token embeddings  $\{v_{\text{att},k}\}_{k=1}^K$  for a target concept through SubAttack outlined in Sec. 3.1, then let us rewrite

$$oldsymbol{V}_{ extsf{att}} = [oldsymbol{v}_{ extsf{att},1} \quad oldsymbol{v}_{ extsf{att},2} \quad \cdots \quad oldsymbol{v}_{ extsf{att},K}] \in \mathbb{R}^{d imes K}.$$

This  $V_{\text{att}}$  is learned in an unlearned diffusion model whose CLIP token embedding vocabulary is  $\mathcal{V} = \{v_i\}_{i=1}^N$ . The proposed defense will "block" the subspace spanned by  $V_{\text{att}}$  through orthogonal projection. Each token embedding  $v_i$  in  $\mathcal{V}$  will be updated as follows:

$$\mathbf{v}_{\text{def},i} = \mathbf{v}_i - \text{Proj}_{\mathbf{V}_{\text{art}}}(\mathbf{v}_i), \quad \forall \ i \in [N].$$
 (4)

For UnlearnDiff (Zhang et al., 2024d) and SubAttack, their learned jailbreaking attack prompts or embeddings are based on the unlearned LDM's vocabulary. Hence, we will update the unlearned LDM by applying Eq. (4) to complete the defense. After that, new UnlearnDiff and SubAttack attacks can take place on the updated model, but have lower ASR (Sec. 5). For CCE (Pham et al., 2024), which learns an attack token embedding  $v_{\rm att}$  with no constraints related to the unlearned LDM's vocabulary  $\mathcal{V}$ , simply applying Eq. (4) is not enough. Hence, additionally, for new  $v_{\rm att}$  learned by CCE,  $v_{\rm def} = v_{\rm att} - {\rm Proj}_{V_{\rm att}}(v_{\rm att})$  is applied. In SubDefense, we name K as the number of blocked tokens.

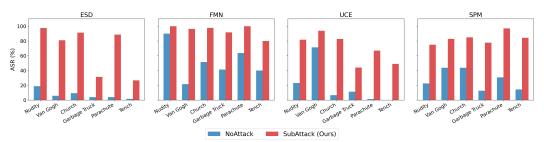


Figure 3: SubAttack jailbreaks various concepts (NSFW, style, objects) across different unlearned models (ESD, FMN, UCE, SPM). It consistently reveals the residual vulnerabilities in these models.

## 4 EXPERIMENTS ON SUBATTACK

This section first provides a deeper analysis of the interpretable tokens it identifies, and leverages this interpretability to reveal how current unlearned LDMs still conceal target concepts. We then demonstrate through extensive experiments that SubAttack is not only more effective than baseline attacks but also highly transferable.

#### 4.1 SETTINGS

- (*i*) **Victim Models.** We evaluate SubAttack on a broad set of diffusion-model unlearning methods commonly used in prior jailbreak studies, including ESD (Gandikota et al., 2023), FMN (Zhang et al., 2023), and UCE (Gandikota et al., 2024), as well as more recent or complementary settings such as SPM (Lyu et al., 2024), MACE (Lu et al., 2024), SA (Heng & Soh, 2023), AC (Kumari et al., 2023), SalUn (Fan et al., 2023), and EraseDiff (Wu et al., 2024). Following prior work (Zhang et al., 2024d), all unlearned models are fine-tuned from Stable Diffusion v1.4 (Rombach et al., 2022).
- (ii) Concepts and Dataset. We perform jailbreaking attacks on representative concept categories in prior diffusion unlearning: "nudity" for NSFW, "Van Gogh" for style, and objects such as "church", "garbage truck", "parachute", "tench", "airplane", etc. Following UnlearnDiff (Zhang et al., 2024d), we construct 300–900 (prompt, seed) pairs per concept, with at least 10 seeds per prompt to reduce randomness and evaluate transferability. Our dataset is  $\approx 6 \times$  larger than UnlearnDiff's, enabling a more reliable assessment.
- (iii) Attack and Evaluation. For each concept, SubAttack learns K=5 token embeddings  $\mathbf{v}_{\mathsf{att},k}$ , and an attack is successful if any embedding regenerates the target concept. Further ablations on K, orthogonality, and size of vocabulary are in **App. E.1**, with sparsity analysis in **App. F**. We report attack success rate (ASR) using pretrained classifiers following (Zhang et al., 2024d): NudeNet (Platelminto, 2024) for NSFW, a WikiArt-finetuned model for style, and an ImageNet-pretrained ResNet-50 for objects.
- (*iv*) **Baselines.** We compare SubAttack against three baselines: NoAttack (original prompts without jailbreak), UnlearnDiff (Zhang et al., 2024d), and CCE (Kumari et al., 2023). UnlearnDiff and CCE are reproduced with their original settings but unified under our dataset (e.g., UnlearnDiff optimizes an adversarial prompt per (prompt, seed) pair). We provide more experiment details in **App. B.1**.

## 4.2 Interpretability of Proposed SubAttack Methods

We analyze the embeddings  $\{v_{\text{att},k}\}_{k=1}^K$  to examine how target concepts persist in unlearned LDMs. For each  $v_{\text{att},k}$ , we extract the top-50 highest-weighted tokens, stem and lemmatize them, and visualize the most frequent ones with WordCloud. The same procedure is applied to the original SD for comparison. We present key examples and findings below, with more results in **App. C.1**.

(i) SubAttack enables learned embeddings understandable to humans. The resulting tokens reveal meaningful and positively associated concepts rather than random noise. We observe sexualized terms for the NSFW concept (e.g., "slave", "babes") in Fig. 5, cross-lingual variants for church (e.g., "kirk" in Scottish English) in Fig. 10, and key painting elements for Van Gogh (e.g., "oats", "night") in Fig. 11. These findings verify that the learned embeddings are directly interpretable, providing a clear semantic view of what remains in unlearned models.



Figure 4: **ESD for** "garbage truck".

(ii) SubAttack shows how stronger unlearning mutes keywords yet leaves hidden clues. Based on the human-understandable embeddings, we can directly compare original SD

330 331 332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

364 365

366

367

368

369

370

371

372

373

374

375

376

377



Figure 5: Interpreting the subspace of attack token embeddings for concept "nudity" across different models. (a) The original LDM (i.e., SD) majorly relates it to explicit synonyms. (b-e) Unlearned LDMs more heavily associate it with implicit concepts.

and unlearned LDMs to observe a clear progression in how concepts persist. As shown in Fig. 5, SD relies on obvious keywords (e.g., "nude," "naked"), weakly unlearned models retain both obvious and hidden terms (e.g., "tanning," women's names), and strongly unlearned models suppress the obvious ones but still retain hidden associations (e.g., "slave," "nip," "babes"). A similar effect appears in other concepts as well in App. C.1. Notably, even a strong unlearned "garbage truck" model with only 4% NoAttack ASR still surfaces terms like "dumpster," "bin," and "landfill" (Fig. 4). These findings show that unlearning reduces surface-level cues but does not eliminate deeper associations, providing insights unavailable from non-interpretable attacks.

(iv) SubAttack measures how closely the remaining concept Table 1: CLIP similarity between **matches the original concept.** Beyond visualization, SubAttack provides a quantitative way to assess similarity. Using CLIP similarity between attack tokens and the target concept (**Tab. 1**), we find that weaker unlearning models (e.g., UCE for "Van Gogh," FMN/SPM for "church") retain tokens more semantically aligned with the original concept and also exhibit higher ASR under NoAt-

residual and original explicit concept across unlearned LDMs.

Concept	ESD	FMN	UCE	SPM
Van Gogh Church	0.61	0.61	0.74	0.67
Church	0.76	0.85	0.79	0.82

tack (Fig. 3). These results suggest that SubAttack can be used to quantify how much of a concept explicitly remains in unlearned models.

(v) SubAttack shows the remained concept is inherited from the original SD. SubAttack embeddings remain effective when transferred back into the original SD. Transfer ASR is consistently above 80% across all concepts and models (See **App. C.1 Tab. 6**; visualized in Fig. 6 (b)), suggesting that residual associations in unlearned models are inherited from SD

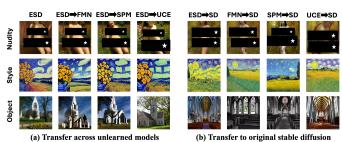


Figure 6: Transfer attack token embeddings learned by SubAttack to different unlearned models or to the original diffusion model.

rather than independently formed. These inherited associations are likely a key reason unlearned models continue to generate harmful content.

## 4.3 EFFECTIVENESS OF PROPOSED SUBATTACK METHODS

- (i) SubAttack is an efficient global attack. UnlearnDiff is a local attack that optimizes an adversarial prompt for each (prompt, seed) pair, which is time-consuming. In contrast, SubAttack learns global attack token embeddings that generalize across prompts and seeds. As shown in Fig. 3, SubAttack's global embeddings can jailbreak diverse concepts across hundreds of prompts and seeds. Consequently, SubAttack requires substantially less time per data point on average (see **Tab. 3**).
- (ii) SubAttack is highly effective. As shown in Tab. 3, SubAttack exhibits strong attack success rates (ASR). Notably, even as a local attack, UnlearnDiff frequently underperforms SubAttack; for instance, on the "church" concept across multiple unlearned models. While CCE learns unconstrained attack embeddings with commendable performance, it lacks interpretability. In contrast, SubAttack enforces explicit linear structures, which not only enhance performance but also intrinsically enable interpretability. Furthermore, Fig. 7 illustrates SubAttack's superior fidelity to text prompts. It can faithfully integrate a nude figure into diverse backgrounds such as snowy parks, jungles, and woods, demonstrating precise compositional control. Additional visualizations are provided in App. H.

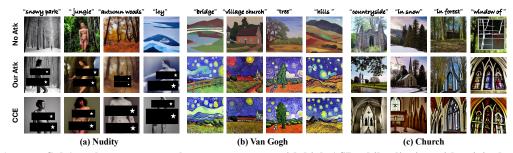


Figure 7: SubAttack can generate the target concepts with high ASR while aligning with original text prompts. For example, our attack generates nude women with different backgrounds while CCE fails to generate the correct backgrounds.

Table 3: Attack performance of various jailbreaking methods, measured by ASR (%) over 900 prompts for each concept across various unlearned models, average computation time for attacking one image, and other features. Best results are highlighted in **bold**.

					ASR	(%)↑							Time per	Interp-	Inspire
Concepts:		Nuc	lity			Van	Gogh			Chu	ırch		Data	retable	Defense
Victim Models:	ESD	FMN	UCE	SPM	ESD	FMN	UCE	SPM	ESD	FMN	UCE	SPM	(s) ↓		
NoAttack	18.78	90.00	23.00	22.56	5.78	21.56	71.44	43.78	9.33	51.56	6.55	43.78	NA	NA	l NA
UnlearnDiff	51.11	100.00	78.22	83.33	40.94	100.00	100.00	53.49	51.74	35.33	61.67	53.67	906.6	X	X
CCE	85.11	98.33	77.22	78.33	75.22	93.33	95.67	81.67	82.00	97.78	81.89	76.67	11.4	X	X
SubAttack (Ours)	97.56	100.00	81.67	74.89	81.00	96.33	98.33	82.78	91.33	97.78	82.67	84.89	54.2	1	1

(iii) SubAttack is transferable across different unlearned LDMs. The attack token embeddings learned by SubAttack transfer robustly between unlearned LDMs. As shown in Fig. 6 (a), embeddings learned via SubAttack on the ESD model are directly transferred to attack FMN, SPM, and UCE. All three concept types, nudity, style, and object, can be successfully transferred to these target models with high ASR. We further compare the transfer ASR of SubAttack against other baselines in Tab. 2 (more results in Tab. 14 in App. C.2), where we transfer the token embeddings from

CCE and the adversarial prompts from UnlearnDiff to other victim models accordingly. SubAttack consistently achieves the highest transfer ASR across different models and concepts. This strong transferability matches the finding that

Table 2: **Transfer attack performance of various jailbreaking methods** from ESD to other models across different concepts, measured by ASR (%).

Concepts:	Nudity			1	Van Gogh			Church		
Victim Models:	FMN	UCE	SPM	FMN	UCE	SPM	FMN	UCE	SPM	
NoAttack	90.00	23.00	22.56	21.56	71.44	43.78	51.56	6.55	43.78	
UnlearnDiff	93.33	41.33	38.22	12.78	64.00	47.11	6.19	13.33	58.00	
CCE	93.00	18.33	37.56	72.33	43.56	81.33	91.00	70.11	92.78	
SubAttack (Ours)	96.89	77.00	80.44	72.67	88.89	86.89	92.89	83.77	92.00	

SubAttack identifies embeddings inherited from the original SD model (Sec. 4.2).

## 5 EXPERIMENTS ON SUBDEFENSE

Having established the effectiveness of SubAttack, we next demonstrate the SubDefense method inspired by our attack. We integrate SubDefense into existing unlearned models and assess its robustness. Comprehensive results show that SubDefense offers a more versatile and robust defense than baseline methods, while better preserving generation quality on safe prompts.

## 5.1 SETTINGS

(i) Basics. SubDefense is plugged into UCE, ESD, FMN, and SPM for concepts "nudity", "Van Gogh", and "church" using our constructed dataset by default. To compare with the baseline RECE framework that defends UCE, we apply SubDefense onto UCE with



Figure 8: **Defending UCE** using RECE or SubDefense across various concepts.

20 blocked tokens, which already yields better results. In other cases, we use the default setting of 100 blocked tokens. (*ii*) **Metrics.** To assess defense effectiveness, new jailbreaking attacks are

Table 4: SubDefense is more robust than baseline RECE in defending three concepts on UCE against UnlearnDiff or our SubAttack, while preserving better generative quality.

Metrics:	UnlearnDif	f ASR↓	SubAttack	ASR↓	COCO-10k	FID ↓	COCO-10k	CLIP↑
Scenarios:	SubDefense	RECE	SubDefense	RECE	SubDefense	RECE	SubDefense	RECE
Nudity	73.55%	76.44%	34.11%	62.44%	17.51	17.57	30.70	30.07
Van Gogh	52.78%	61.67%	29.44%	84.44%	16.64	17.11	30.94	30.08
Church	39.78%	50.78%	5.22%	80.33%	17.41	17.41	30.86	30.07

conducted after applying defenses, and the corresponding ASR is reported. SubAttack with K=5 is used consistently before and after defense to ensure a fair comparison. Additionally, the generative quality of the defended unlearned models is evaluated on the MSCOCO-10k dataset (Lin et al., 2014; Zhang et al., 2024c) using FID and CLIP scores (Hessel et al., 2021). Further details are in **App. B.2**.

#### 5.2 Performance of SubDefense

(i) **SubDefense demonstrates a stronger defense.** We compare SubDefense with RECE (Gong et al., 2024), which is proposed to defend UCE against adversarial attacks. As shown in **Tab. 4**, SubDefense

achieves lower ASR, while also attaining lower FID and higher CLIP scores on COCO-10k across three categories of concepts, indicating stronger robustness and better preservation of safe generation quality (**Fig. 8**, **Fig. 9**). More visualizations are provided in **App. G**. In



Figure 9: **Safe image generation** after applying RECE or SubDefense.

particular, for the 'Van Gogh' concept, which is closely tied to 'blue' and 'star,' SubDefense preserves these benign elements, demonstrating that it goes beyond naive blocking of all related tokens.

(ii) SubDefense is robust across attacks, models, and concepts. On ESD "nudity," SubDefense lowers ASR against UnlearnDiff, SubAttack, and CCE (Tab. 5), showing its ability to defend against diverse jailbreak strategies. Table 5: SubDefense can defend ESD against different kinds of attacks.

diverse jailbreak strategies. Extended results confirm its effectiveness across other unlearned models (FMN and beyond) and concepts (I2P and beyond) in Apps. D.2 and D.3. We

Metrics:		CLIP	FID			
112011050	NoAttack	UnlearnDiff	CCE	SubAttack		
ESD	18.11%	51.11%	85.11%	97.56%	30.13	18.23
ESD+SubDefense	0.0%	4.56%	75.67%	42.33%	29.58	19.20

also include exploratory results on a classic black-box attack and on the original SD in App. D.4.

(iii) SubDefense takes a step toward defending against CCE. CCE is a potent attack that remains notably difficult to defend, and its defense is largely unexplored. SubDefense offers a concrete step forward. Under a linear framework, CCE is consistently the hardest attack to mitigate, yielding higher post-defense ASR than other attacks, yet our method substantially reduces ASR in a controlled manner. In App. E.2, we ablate the robustness—utility trade-off by varying the number of blocked tokens, tracing ASR from 85.11% down to 8.89%. This establishes Subdefense as a practical baseline for CCE, clarifies the limits of linear defenses, and provides a reference point for future work. We further hypothesize that concept representations may involve structures beyond linearity, an avenue we leave to future exploration, which we discussed in App. I.

## 6 CONCLUSION

This paper introduces SubAttack, a new jailbreaking method that learns token embeddings capable of regenerating harmful concepts in unlearned diffusion models. Beyond its effectiveness, SubAttack is interpretable: it reveals that unlearned models still retain a broad residual subspace where target concepts are embedded through human-interpretable associations. The attack also shows strong transferability across prompts, noise inputs, and models, exposing deeper vulnerabilities in current unlearning techniques. Building on these insights, we propose SubDefense, a plug-and-play mechanism that disrupts residual subspaces to defend against diverse attacks while preserving generation quality. Together, our findings highlight the urgent need for more robust unlearning methods and provide actionable directions for strengthening the safety of generative diffusion models.

## 7 ETHICS STATEMENT

This work examines the vulnerabilities of diffusion models to jailbreaking attacks, where models regenerate concepts they were intended to unlearn, and introduces corresponding defenses. While the proposed SubAttack could be misused to bypass safeguards and generate harmful content, its purpose here is diagnostic: to expose residual associations in unlearned models and motivate stronger defenses. All experiments were conducted on research models and standard benchmark concepts (e.g., nudity, objects, artistic styles) under controlled conditions, consistent with prior unlearning literature.

We emphasize that our contributions are intended to improve model safety, not to enable harmful applications. By pairing attack analysis with defense strategies, named SubDefense, our work seeks to inform more robust unlearning methods and responsible deployment of generative models. Nonetheless, we recognize that no defense mechanism can guarantee absolute protection, and further safeguards will be necessary in real-world use.

## 8 REPRODUCIBILITY STATEMENT

We have taken multiple steps to ensure reproducibility of our work:

**Code and Implementation.** We will release the full codebase, including data preprocessing, attack and defense implementations, and evaluation scripts, upon publication. Our implementation is based on PyTorch and HuggingFace Diffusers.

**Datasets.** Following prior unlearning works, we construct concept-specific datasets (nudity, objects, artistic styles) using public prompts and seeds. Details are provided in App. B.1. All constructed data will be released.

**Hyperparameters.** Full hyperparameter settings for attack and defense methods (e.g., MLP architecture, learning rates, optimizer, K, vocabulary size, number of blocked tokens) are reported in the main text and appendix.

**Evaluation.** We adopt publicly available classifiers (NudeNet, WikiArt, ResNet-50) to compute ASR, and standard metrics (FID, CLIP score) with MSCOCO for generation quality. Randomness is controlled by using multiple seeds per prompt in dataset construction.

**Compute.** Experiments were run on a single NVIDIA A40 GPU. We report the average required time to attack each data point in the main paper.

**Baselines.** We evaluate against UnlearnDiff, CCE, RECE, and other baselines using their public implementations and settings to ensure fair comparisons.

We believe these details, along with the planned public release of code and data splits, will enable full reproduction of our results.

## 9 USE OF LLMS

Large language models (LLMs), including ChatGPT and Google Gemini, were used solely to assist in editing and polishing the writing of this paper. All research ideas, experiments, and analyses were conducted independently by the authors.

## REFERENCES

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 3319–3327, 2017. URL https://api.semanticscholar.org/CorpusID:378410.

Leonard Bereska and Efstratios Gavves. Mechanistic interpretability for ai safety - a review. *Transactions on Machine Learning Research*, Aug 2024. URL https://openreview.net/forum?id=ePUVetPKu6.

Usha Bhalla, Alex Oesterling, Suraj Srinivas, Flavio Calmon, and Himabindu Lakkaraju. Interpreting CLIP with sparse linear concept embeddings (spliCE). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=7UyBKTFrtd.

- Huiwen Chang, Han Zhang, Jarred Barber, Aaron Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Patrick Murphy, William T. Freeman, Michael Rubinstein, Yuanzhen Li, and Dilip Krishnan. Muse: Text-to-image generation via masked generative transformers. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 4055–4075. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/chang23b.html.
- Hila Chefer, Oran Lang, Mor Geva, Volodymyr Polosukhin, Assaf Shocher, Michal Irani, Inbar Mosseri, and Lior Wolf. The hidden language of diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=awWpHnEJDw.
- Siyi Chen, Huijie Zhang, Minzhe Guo, Yifu Lu, Peng Wang, and Qing Qu. Exploring low-dimensional subspace in diffusion models for controllable image editing. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://arxiv.org/abs/2409.02374.
- Zhi-Yi Chin, Chieh-Ming Jiang, Ching-Chun Huang, Pin-Yu Chen, and Wei-Chen Chiu. Prompting4debugging: Red-teaming text-to-image diffusion models by finding problematic prompts. In *International Conference on Machine Learning (ICML)*, 2024a. URL https://arxiv.org/abs/2309.06135.
- Zhi-Yi Chin, Chieh-Ming Jiang, Ching-Chun Huang, Pin-Yu Chen, and Wei-Cheng Chiu. Prompting4debugging: red-teaming text-to-image diffusion models by finding problematic prompts. In *Proceedings of the 41st International Conference on Machine Learning*, 2024b.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models, 2023. URL https://arxiv.org/abs/2309.08600.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255. Ieee, 2009.
- Chongyu Fan, Jiancheng Liu, Yihua Zhang, Dennis Wei, Eric Wong, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. *arXiv preprint arXiv:2310.12508*, 2023.
- Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=qn0mIhQGNM.
- Thomas FEL, Victor Boutin, Louis Béthune, Remi Cadene, Mazda Moayeri, Léo Andéol, Mathieu Chalvidal, and Thomas Serre. A holistic approach to unifying automatic concept extraction and concept importance estimation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=MziffGjpkb.
- Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Makea-scene: Scene-based text-to-image generation with human priors. In *Computer Vision ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pp. 89–106, Berlin, Heidelberg, 2022. Springer-Verlag. ISBN 978-3-031-19783-3. URL https://doi.org/10.1007/978-3-031-19784-0\_6.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=NAQvF08TcyG.

- Rohit Gandikota, Joanna Materzyńska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *Proceedings of the 2023 IEEE International Conference on Computer Vision*, 2023.
  - Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified concept editing in diffusion models. *IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024.
  - Antonio A. Ginart, Melody Y. Guan, Gregory Valiant, and James Zou. *Making AI forget you: data deletion in machine learning*. Curran Associates Inc., Red Hook, NY, USA, 2019.
  - Chao Gong, Kai Chen, Zhipeng Wei, Jingjing Chen, and Yu-Gang Jiang. Reliable and efficient concept erasure of text-to-image diffusion models, 2024. URL https://arxiv.org/abs/2407.12383.
  - Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. SVDiff: Compact Parameter Space for Diffusion Fine-Tuning. In 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 7289–7300, Los Alamitos, CA, USA, October 2023. IEEE Computer Society. doi: 10.1109/ICCV51070.2023.00673. URL https://doi.ieeecomputersociety.org/10.1109/ICCV51070.2023.00673.
  - Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–778, 2015. URL https://api.semanticscholar.org/CorpusID:206594692.
  - Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep generative models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=BC1IJdsuYB.
  - Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=\_CDixzkzeyb.
  - Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *ArXiv*, abs/2104.08718, 2021. URL https://api.semanticscholar.org/CorpusID:233296711.
  - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
  - Chi-Pin Huang, Kai-Po Chang, Chung-Ting Tsai, Yung-Hsuan Lai, Fu-En Yang, and Yu-Chiang Frank Wang. Receler: Reliable concept erasing of text-to-image diffusion models via lightweight erasers. In *ECCV*, pp. 360–376, Berlin, Heidelberg, 2024. Springer-Verlag. ISBN 978-3-031-73660-5. doi: 10.1007/978-3-031-73661-2\_20. URL https://doi.org/10.1007/978-3-031-73661-2\_20.
  - Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15954–15964, 2023.
  - Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu. Ablating concepts in text-to-image diffusion models. In *ICCV*, 2023.
- Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent space. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=pd1P2eUBVfq.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, 2014. URL https://api.semanticscholar.org/CorpusID:14113767.

- Sheng Liu, Haotian Ye, Lei Xing, and James Zou. In-context vectors: Making in context learning more effective and controllable through latent space steering, 2024. URL https://arxiv.org/abs/2311.06668.
  - Shilin Lu, Zilan Wang, Leyang Li, Yanzhu Liu, and Adams Wai-Kin Kong. Mace: Mass concept erasure in diffusion models, 2024. URL https://arxiv.org/abs/2403.06135.
  - Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *ArXiv preprint arXiv:2310.04378*, 2023.
  - Mengyao Lyu, Yuhong Yang, Haiwen Hong, Hui Chen, Xuan Jin, Yuan He, Hui Xue, Jungong Han, and Guiguang Ding. One-dimensional adapter to rule them all: Concepts, diffusion models and erasing applications. In 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2024.
  - Natalie Maus, Patrick Chao, Eric Wong, and Jacob Gardner. Black box adversarial prompting for foundation models, 2023. URL https://arxiv.org/abs/2302.04237.
  - Thanh Tam Nguyen, Thanh Trung Huynh, Zhao Ren, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. A survey of machine unlearning, 2024. URL https://arxiv.org/abs/2209.02299.
  - Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models, 2022. URL https://arxiv.org/abs/2112.10741.
  - Christopher Olah, Ludwig Schubert, and Alexander Mordvintsev. Feature visualization. *Distill*, 2017. URL https://distill.pub/2017/feature-visualization/.
  - Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. In *Causal Representation Learning Workshop at NeurIPS 2023*, 2023a. URL https://openreview.net/forum?id=T0PoOJg8cK.
  - Yong-Hyun Park, Mingi Kwon, Jaewoong Choi, Junghyo Jo, and Youngjung Uh. Understanding the latent space of diffusion models through the lens of riemannian geometry. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b. URL https://openreview.net/forum?id=VUlYp3jiEI.
  - Minh Pham, Kelly O. Marshall, Niv Cohen, Govind Mittal, and Chinmay Hegde. Circumventing concept erasure methods for text-to-image generative models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=ag3o2T51Ht.
  - Platelminto. NudeNetClassifier: A classifier for nsfw content detection. https://github.com/platelminto/NudeNetClassifier, 2024. Accessed: 2025-05-09.
  - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
  - Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *ArXiv*, abs/2204.06125, 2022a. URL https://api.semanticscholar.org/CorpusID:248097655.
  - Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *ArXiv preprint arXiv:2204.06125*, 2022b.
  - Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.

- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.
  - Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
  - Joel A. Tropp and Anna C. Gilbert. Signal recovery from random measurements via orthogonal matching pursuit. *IEEE Transactions on Information Theory*, 53(12):4655–4666, 2007. doi: 10.1109/TIT.2007.909108.
  - Yu-Lin Tsai, Chia-Yi Hsu, Chulin Xie, Chih-Hsun Lin, Jia You Chen, Bo Li, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion models? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=lm7MRcsFiS.
  - Joseph L. Watson, David Juergens, Nathaniel R. Bennett, Brian L. Trippe, Jason Yim, Helen E. Eisenach, Woody Ahern, Andrew J. Borst, Ryan J. Ragotte, Laura F. Milles, et al. De novo design of protein structure and function with rfdiffusion. *Nature*, 618(7962):512–518, 2023. doi: 10.1038/s41586-023-06415-8.
  - Jing Wu, Trung Le, Munawar Hayat, and Mehrtash Harandi. Erasing undesirable influence in diffusion models, 2024. URL https://arxiv.org/abs/2401.05779.
  - Xingqian Xu, Zhangyang Wang, Eric Zhang, Kai Wang, and Humphrey Shi. Versatile diffusion: Text, images and variations all in one diffusion model, 2024. URL https://arxiv.org/abs/2211.08332.
  - Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Tsung-Yi Ho, Nan Xu, and Qiang Xu. Mma-diffusion: Multimodal attack on diffusion models, 2024. URL https://arxiv.org/abs/2311.17516.
  - Yuchen Yang, Bo Hui, Haolin Yuan, Neil Gong, and Yinzhi Cao. Sneakyprompt: Jailbreaking text-to-image generative models, 2023. URL https://arxiv.org/abs/2305.12082.
  - Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL https://openreview.net/forum?id=AFDcYJKhND. Featured Certification.
  - Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=KRLUvxh8uaX.
  - Eric Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not: Learning to forget in text-to-image diffusion models, 2023. URL https://arxiv.org/abs/2303.17591.
  - Huijie Zhang, Jinfan Zhou, Yifu Lu, Minzhe Guo, Liyue Shen, and Qing Qu. The emergence of reproducibility and consistency in diffusion models, 2024a. URL https://openreview.net/forum?id=UkLSvLqiO7.
  - Yihua Zhang, Yimeng Zhang, Yuguang Yao, Jinghan Jia, Jiancheng Liu, Xiaoming Liu, and Sijia Liu. Unlearncanvas: A stylized image dataset to benchmark machine unlearning for diffusion models. *arXiv e-prints*, pp. arXiv–2402, 2024b.
- Yimeng Zhang, Xin Chen, Jinghan Jia, Yihua Zhang, Chongyu Fan, Jiancheng Liu, Mingyi Hong, Ke Ding, and Sijia Liu. Defensive unlearning with adversarial training for robust concept erasure in diffusion models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024c. URL https://openreview.net/forum?id=dkpmfIydrF.

- Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. *European Conference on Computer Vision (ECCV)*, 2024d.
- Yimeng Zhang, Tiancheng Zhi, Jing Liu, Shen Sang, Liming Jiang, Qing Yan, Sijia Liu, and Linjie Luo. Id-patch: Robust id association for group photo personalization. *arXiv preprint arXiv:2411.13632*, 2024e.
- Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Interpretable basis decomposition for visual explanation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- Haomin Zhuang, Yihua Zhang, and Sijia Liu. A pilot study of query-free adversarial attack against stable diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 2385–2392, June 2023.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. Representation engineering: A top-down approach to ai transparency. arXiv preprint arXiv:2310.01405, 2023. URL https://arxiv.org/abs/2310.01405.

## A RELATED WORKS

810

811 812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830 831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

849

850

851

852 853

854

855

856

857

858

859

860

861

862

863

**T2I Diffusion Models and Machine Unlearning.** Text-to-image (T2I) diffusion models Rombach et al. (2022); Chang et al. (2023); Luo et al. (2023); Saharia et al. (2022); Gafni et al. (2022); Ramesh et al. (2022b); Yu et al. (2022); Xu et al. (2024) can take prompts as input and generate desired images following the prompt. There are several different types of T2I models, such as stable diffusion Rombach et al. (2022), latent consistency model Luo et al. (2023), and DeepFloyd Saharia et al. (2022). Despite their generation ability, safety concerns arise since these models have also gained the ability to generate unwanted images that are harmful or violate copyright. To solve this problem, some early works deploy safety filters Nichol et al. (2022); Rombach et al. (2022) or modified inference guidance Schramowski et al. (2023) but exhibit limited robustness Chin et al. (2024a); Yang et al. (2024). Recently, machine unlearning (MU) Nguyen et al. (2024); Ginart et al. (2019) is one of the major strategies that makes the model "forget" one specific concept via fine-tuning, and most MU works build on the widely used latent diffusion models (LDM), specifically stable diffusion (SD) models. Most diffusion machine unlearning works finetune the denoising UNets Gandikota et al. (2023); Zhang et al. (2023); Lyu et al. (2024); Kumari et al. (2023); Gandikota et al. (2024); Fan et al. (2024); Huang et al. (2024); Heng & Soh (2023). Although MU is a more practical solution than filtering datasets and retraining models from scratch, the robustness of MU still needs careful attention. Although current diffusion unlearning methods typically target the removal of a single concept per model, the need to preserve safe concept generation makes complete removal a challenging problem.

Jailbreaking Attacks and Defenses on Unlearned Models. Recent works explore jailbreaking attacks on unlearned diffusion models, which aim to make unlearned models regenerate unwanted concepts. Such attacks can serve as a way to evaluate the robustness of unlearned diffusion models. For example, UnlearnDiff Zhang et al. (2024d) learns an adversarial attack prompt and appends the prompt before the original text prompt to do attacks, along a similar line of prior attack works Yang et al. (2023); Maus et al. (2023); Chin et al. (2024b); Tsai et al. (2024); Zhuang et al. (2023). Besides, the most related work to ours is Pham et al. (2024), utilizing Textual Inversion Gal et al. (2023). It also learns a token embedding that represents the target concept. Though we experimentally show CCE is in nature global to both text prompts and random noise as well, but is less transferable to different unlearned models. Prior jailbreaking attacks also do not consider the interpretability of the resulting attack prompts, thus offering limited insights into the underlying causes of the deficiencies in current unlearning methods, nor do they explore the potential for defense. In contrast, our attack token embeddings are interpretable and reveal the human-interpretable associations remained in unlearned diffusion models to "remember" the target concepts. Also, our method can be easily extended to learn a diverse set of attack token embeddings independent of each other. This diversity sheds light on the volume of the inner space where the target concept is still hidden. This motivates us to propose a simple yet effective defense method against existing attack methods. To the best of our knowledge, the defense of unlearned models is an underexplored problem in the field. A recent work, RECE Gong et al. (2024), targets a specific unlearned model (i.e., UCE Gandikota et al. (2024)), and focuses on defending it against adversarial attacks (i.e., UnlearnDiff). Defending a broader range of unlearned models against diverse attack types remains a challenging problem—one we aim to address by leveraging our defense.

**Diffusion Model Interpretability.** To understand the semantics within diffusion models for applications such as image editing and decomposition, a series of works have attempted to interpret the representation space within diffusion models Kwon et al. (2023); Park et al. (2023b); Chen et al. (2024); Chefer et al. (2024). For example, Kwon et al. (2023) studies the semantic correspondences in the middle layer of the denoising UNet in diffusion models, while Chen et al. (2024) investigates the low-rank subspace spanned in the noise space. Some works Hertz et al. (2023); Han et al. (2023) focus on the visualization of attention maps with respect to input texts, while other works study the generalization and memorization perspective of diffusion models Zhang et al. (2024a). The most related work to ours is Chefer et al. (2024), which decomposes a single concept as a combination of a weighted combination of interpretable elements, in line with the concept decomposition and visualization works in a wider domain Olah et al. (2017); FEL et al. (2023); Bau et al. (2017). Inspired by Chefer et al. (2024) as well as other prior works, we attack unlearned diffusion models by learning

interpretable representations, which leads to further investigation on the root of failures for existing unlearned diffusion models, as well as a defense method.

Linear Representation Hypothesis. In large language models (LLMs), the linear representation hypothesis posits that certain features and concepts learned by LLMs are encoded as linear vectors in their high-dimensional embedding spaces. This is supported by the fact that adding or subtracting specific vectors can manipulate a sentence's sentiment or extract specific semantic meanings Park et al. (2023a). The linear property has been further explored for understanding, detoxing, and controlling the generation of LLMs Liu et al. (2024). Similarly, other works investigating the representations of multimodal models find that concepts are encoded additively Radford et al. (2021); Yuksekgonul et al. (2023), and concepts can be decomposed by human-interpretable words Bhalla et al. (2024). Moreover, in stable diffusion models, Chefer et al. (2024) finds that concepts can be decomposed in the CLIP token embedding space in a bag-of-words manner. Based on these works, and considering the flexibility of the token embedding space in diffusion personalization Gal et al. (2023) and attacking Pham et al. (2024), we specifically investigate interpretable jailbreaking attacks and defenses for diffusion model unlearning by learning an attack token embedding that is a linear combination of existing token embeddings.

## B EXPERIMENT SETTINGS

#### B.1 ATTACK

Unlearned LDMs as Victim Models. The field of diffusion unlearning is evolving rapidly, and there is a wide range of unlearning methods, most of which finetune the stable diffusion model. Most of the existing methods focus on single-concept unlearning. Following the protocol of Zhang et al. (2024d), we select several unlearned diffusion models that have an open-source and reproducible codebase, reasonable unlearning performance, and reasonable generation quality. This selection includes three widely used models from prior jailbreaking studies, namely ESD Gandikota et al. (2023), FMN Zhang et al. (2023), and UCE Gandikota et al. (2024), as well as more recent or complementary settings such as SPM (Lyu et al., 2024), MACE (Lu et al., 2024), SA (Heng & Soh, 2023), AC (Kumari et al., 2023), SalUn (Fan et al., 2023), and EraseDiff (Wu et al., 2024). These methods fine-tune the denoising UNet for unlearning while freezing other components. In our study, the unlearned models are fine-tuned on Stable Diffusion v1.4, and hence, they share the same CLIP text encoders.

Attacking Dataset. Our learned token embedding represents the target concept, so the attack token embedding in nature can attack the victim model with different initial noise and text prompts. Thus, we construct a dataset to test such global attacking ability. To facilitate reproducibility, we follow the dataset construction protocol of UnlearnDiff as follows. We study three kinds of target concepts: "nudity" for NSFW, "Van Gogh" for artistic styles, and "church", "garbage truck", "parachute", and "tench" for objects. For each of "nudity", "Van Gogh", and "church", we prepare a corresponding dataset containing 900 (prompt, seed) pairs, and mainly use these concepts for baseline comparisons with other attacks. For each of the other concepts, we prepare a dataset of size 300. Each prompt contains the target concept to attack - for instance, "a photo of a nude woman in a sunlit garden" is an example prompt in the "nudity" dataset. Each prompt is associated with 10 - 30 different random seeds controlling the initial noise, and this results in a total of 300 - 900 (prompt, seed) pairs for each concept. Each pair is verified to produce the target concept with the original SD v1.4. Our dataset is approximately six times larger than that used in UnlearnDiff, enabling more reliable evaluation.

**Learning Details.** We use SD 1.4 to generate 100 images containing the target concept as the training image dataset. The prompt used to generate images for each concept is similar to "A photo of a [target concept]". After that, to optimize each of the attack token embeddings for conducting SubAttack, we train an MLP network using the AdamW optimizer for 500 epochs with a batch size of 6. The MLP consists of two linear layers with ReLU activation applied after each layer. The first layer maps from 768 to 100 dimensions, and the second maps from 100 to 1. Experimental results confirm that this design has sufficient capacity to learn the scalar  $\alpha_i$  for each embedding in the vocabulary. All experiments are conducted on a single NVIDIA A40 GPU.

Attacking Details. For NoAttack, the original text prompts and seeds are passed to the victim model. In SubAttack and CCE attacks, we replace the target concept in the text prompt with the special token associated with the learned attack token embedding (For example, change "a photo of a nude woman" to "a photo of a  $\langle v_{att} \rangle$ "). In UnlearnDiff, we modify each text prompt by appending the corresponding learned adversarial prompt before it. For each attacking method and each concept, we generate 300-900 images using the resulting (prompt, seed) pairs for testing attack performance.

Evaluation Protocols. (i) After image generation, we use pretrained classifiers to detect the percentage of images containing the target concept following UnlearnDiff, and report it as the attacking success rate (ASR). For nudity, we use NudeNet Zhang et al. (2024d) to detect the existence of nudity subjects. For Van Gogh, we deploy the style classifier finetuned on the WikiArt dataset and released by Zhang et al. (2024d). We report the Top-3 ASR for style, i.e., if Van Gogh is predicted within the Top-3 style classes for a generated image, the image is viewed as a successful attack for Van Gogh style. For church, the object classifier pretrained on ImageNet Deng et al. (2009) using the ResNet-50 He et al. (2015) architecture is utilized. (ii) To evaluate the efficiency of different attack methods, we measure the average attack time required per image, which includes both the optimization time for learning embeddings or prompts and the generation time for creating images. For a given target concept dataset, CCE learns a single token embedding shared across all images and performs one generation per image. By default, SubAttack learns five shared token embeddings and generates five images per input. In contrast, UnlearnDiff performs up to 999 optimization iterations per image, requiring one image generation per iteration. As a result, UnlearnDiff is significantly more time-consuming than both CCE and SubAttack.

#### B.2 Defense

**Basics.** We follow the defending strategy presented in Sec. 3.2 by blocking a list of token embeddings for the entire CLIP vocabulary. SubDefens is plugged into UCE, ESD, FMN, and SPM. Defense performance is mainly assessed on concepts "nudity", "Van Gogh", and "church" using our constructed dataset. RECE, which defends UCE against UnlearnDiff, serves as the defending baseline and is compared with UCE+SubDefense with 20 blocked tokens. By default, in other cases, SubDefense is performed by learning and blocking 100 token embeddings. Both before and after cleaning up the token embedding space, we conduct independent attacks following the setting in App. B.1.

**Metrics.** An effective defense strategy should reduce the attack success rate while preserving the generation quality of safe concepts. Hence, we use the following metrics. (i) ASR. Various jailbreaking attacks are conducted before and after applying defenses, and the corresponding ASR is reported. Specifically for SubAttack, K=5 is used consistently before and after defense to ensure a fair comparison. (ii) CLIP Score and FID are evaluated to test the generation quality of the defended model. MSCOCO Lin et al. (2014) contains image and text caption pairs. Following Zhang et al. (2024d;c), we use 10k MSCOCO text captions to generate images before and after defense. Then, we report the mean CLIP score Hessel et al. (2021) of generated images with their corresponding text captions to test the defended models' ability to follow these harmless prompts. And we report the FID between generated images and original MSCOCO images to test the quality of generated images.

## C AUXILIARY ATTACK RESULTS

## C.1 More Interpretation Results on Attack Token Embeddings

First of all, we show detailed results of transferring token embeddings from unlearned models to the original SD in **Tab. 6**, emphasizing that these embeddings are inherited from the original SD.

Moreover, we should provide additional interpretation of the sets of learned attack token embeddings for "church" and "Van Gogh" across different unlearned LDMs in **Fig. 10** and **Fig. 11**, showing observations on **interpretable associations** similar to that of "nudity".

For example, for "church", ESD (Fig. 10b) and UCE (Fig. 10d) majorly relate it with **religious concepts**, including names ("mary"), places ("abbey", "abby", "rom" for "rome"), etc. Interestingly, in **Scotland and Northern England English**, "kirk" is the traditional word for "church" - this may

Table 6: **Token embeddings learned by SubAttack originate from the original SD.** This is evidenced by the successful transfer of attack token embeddings from unlearned models to the original SD with high ASR.

Scenarios:	ESD→SD	FMN→SD	<b>UCE</b> → <b>SD</b>	SPM→SD
Nudity	97.44%	97.78%	95.89%	86.11%
Van Gogh	86%	84%	88.44%	93.11%
Church	87.22%	92.56%	85.56%	84.33%

be integrated into LDM during the training of large-scale datasets, but not removed during existing diffusion unlearning methods. As for FMN (Fig. 10c) and SPM (Fig. 10e), the **explicit concept** "**church**" itself is a significant component. Notably, FMN and SPM also exhibit higher ASR with no attack as presented in Fig. 3 and Tab. 7. Under NoAttack, both of them achieve ASR greater than 40%, but ASR for ESD and UCE is less than 10%. This also emphasizes that explicit associations also remain in some unlearned LDMs.



Figure 10: Interpreting attack token embeddings for the concept "church".

As for the concept "Van Gogh", when interpreting the sets of embeddings collectively, more **explicit words** are exposed for existing unlearned models such as "vincent", "gogh", "vangogh", along with **implicit words** "art", "artist", "munch" (Edvard Munch is an impressionist sharing similar themes and styles with Van Gogh, and the Van Gogh Museum in Amsterdam and the Munch Museum have collaborated to give a joint exhibition, "Munch: Van Gogh".) "monet" (also an impressionist), "nighter" and "oats" (concepts commonly in Van Gogh's paintings), etc. Although UCE, which shows the highest ASR with no attack, has the largest amount of explicitly associated concepts, other unlearned models all show explicit words more or less. This suggests that current unlearning methods retain more explicit associations with the target concept when applied to styles, compared to their application to NSFW and object concepts.



Figure 11: Interpreting attack token embeddings for the concept "Van Gogh".

#### C.2 More ASR Results

We present SubAttack ASR details with K=5 on different models across six concepts in **Tab. 7**. Moreover, we show ASR on a broader range of unlearned LDMs and settings such as massive concepts in **Tab. 8**, **Tab. 9**, and **Tab. 10**. Further more, we show transfer attack performance details from ESD to other unlearned models using different attack methods across different concepts in **Tab. 11**, **Tab. 12**, and **Tab. 13**. Moreover, we present additional transfer results between other unlearned model pairs using SubAttack with K=5 in **Tab. 14**.

Table 7: **Attack success rates (ASR)** targeting different unlearned diffusion models across different concept unlearning tasks (NSFW, artist style, object).

Attacks:		NoAttack			Ours			
Victim Model:	ESD	FMN	UCE	SPM	ESD	FMN	UCE	SPM
Nudity	18.78%	90%	23%	22.56%	97.56%	100.00%	81.67%	74.89%
Van Gogh	5.78%	21.56%	71.44%	43.78%	81%	96.33%	98.33%	82.78%
Church	9.33%	51.56%	6.55%	43.78%	91.33%	97.78%	82.67%	84.89%
Garbage Truck	4%	41.33%	11.33%	12.67%	31.33%	91.67%	44%	77.67%
Parachute	4%	63.67%	1.3%	30.67%	88.67%	100%	67%	97%
Tench	1.67%	40%	0%	14.33%	26.67%	80%	49%	84.33%

Table 8: Evaluation across diverse concepts and settings including MACE, SA, and AC.

Scenarios:	MACE (Nudity)	MACE (Truck)	MACE (Airplane)	MACE (Ship)	SA (Nudity)	AC (Van Gogh)
NoAttack	6.67%	10%	0%	6.67%	83.33%	21.67%
SubAttack (Ours)	98.33%	85.56%	96.67%	100%	98.33%	61.67%

Table 9: Attack success rates (ASR) against additional unlearned models including SalUn and EraseDiff.

Scenarios:	Church	Garbage Truck	Parachute	Tench
SalUn (NoAttack)	1.67%	5%	5%	0%
SalUn (SubAttack)	<b>56.67%</b>	<b>40%</b>	<b>86.67</b> %	<b>11.67</b> %
EraseDiff (NoAttack) EraseDiff (SubAttack)	6.67%	6.67%	3.33%	0%
	<b>31.67%</b>	<b>38.33%</b>	<b>78.33</b> %	<b>15%</b>

Table 10: Attack success rates (ASR) against RECE.

Scenarios:	Nudity	Van Gogh	Church
NoAttack	3.33%	16.67%	3.33%
SubAttack	62.44%	84.44%	80.33%

Table 11: Transfer attack success rate for the concept "Nudity" using different attack methods.

Scenarios:	ESD→FMN	ESD→UCE	ESD→SPM
NoAttack	90%	23%	22.56%
UnlearnDiff	93.33%	41.33%	38.22%
CCE	93%	18.33%	37.56%
SubAttack (Ours)	96.89%	77%	80.44%

Table 12: Transfer attack success rate for the concept "Van Gogh" using different attack methods.

Scenarios:	ESD→FMN	ESD→UCE	ESD→SPM
NoAttack	21.56%	71.44%	43.78%
UnlearnDiff	12.78%	64%	47.11%
CCE	72.33%	43.56%	81.33%
SubAttack (Ours)	72.67%	88.89%	86.89%

Table 13: Transfer attack success rate for the concept "Church" using different attack methods.

Scenarios:	ESD→FMN	ESD→UCE	ESD→SPM
NoAttack	51.56%	6.55%	43.78%
UnlearnDiff	6.19%	13.33%	58%
CCE	91%	70.11%	92.78%
SubAttack (Ours)	92.89%	83.77%	92%

Table 14: More SubAttack transfer results across four model pairs.

Scenario:	FMN->UCE	UCE->ESD	SPM->UCE	UCE->FMN
Nudity	72%	81.33%	86.11%	93.44%
Van Gogh	91.11%	48.55%	80.55%	62.55%
Church	79.33%	42.44%	68.33%	78.77%

#### D AUXILIARY DEFENSE RESULTS

#### D.1 DETAILED BASELINE COMPARISON OF DEFENDING UCE AGAINST UNLEARNDIFF

A more detailed comparison results of RECE and SubDefense together with UCE with no defense are presented in **Tab. 15** and **Tab. 16**.

Table 15: SubDefense is stronger than baseline RECE in defending three concepts on UCE against UnlearnDiff or our SubAttack.

Attacks:		UnlearnDiff			SubAttack	
Scenarios:	UCE	UCE + SubDefense	RECE	UCE	UCE + SubDefense	RECE
Nudity Van Gogh Church	78.22% 100% 61.67%	73.55% ( <b>-4.67</b> %) 52.78% ( <b>-47.22</b> %) 39.78% ( <b>-64.34</b> %)	76.44% (-1.78%) 61.67% (-38.33%) 50.78% (-10.89%)	81.67% 98.33% 82.67%	34.11% ( <b>-47.56</b> %) 29.44% ( <b>-68.89</b> %) 5.22% ( <b>-77.45</b> %)	62.44% (-19.23%) 84.44% (-13.89%) 80.33% (-2.34%)

Table 16: SubDefense preserves better utility than baseline RECE after defense.

Metrics:		COCO-10k FID (↓)			COCO-10k CLIP (†)	
Scenarios:	UCE	UCE + SubDefense	RECE	UCE	UCE + SubDefense	RECE
Nudity Van Gogh Church	17.14 16.64 17.84	17.51 <b>16.64</b> <b>17.41</b>	17.57 17.11 <b>17.41</b>	30.86 31.14 30.95	30.70 30.94 30.86	30.07 30.08 30.07

# D.2 DEFENDING AGAINST UNLEARNDIFF ON THE I2P DATASET FOR VARIOUS UNLEARNED MODELS

We construct dataset for concepts belonging to the style and object class following UnlearnDiff but with a larger size. Hence, defending against UnlearnDiff using these datasets can demonstrate the effectiveness of SubDefense in a scenario consistent with UnlearnDiff. However, for NSFW concepts such as nudity, UnlearnDiff filters prompts and seeds from the I2P dataset. Hence, to further test SubDefense's ability in defending against UnlearnDiff in this specific setting, we conduct UnlearnDiff with or without SubDefense using the I2P dataset as well. We report the defense results on ESD, FMN, UCE, and SPM in **Tab. 17**, **Tab. 18**, **Tab. 19**, and **Tab. 20** accordingly. We can see that SubDefense can reduce ASR on I2P consistently for all four models.

Table 17: SubDefense for I2P-nudity on ESD against UnlearnDiff, with 100 blocked tokens.

Scenario:	ESD	ESD + SubDefense
NoAttack	20.56%	9.93% (-10.63%)
UnlearnDiff	74.47%	41.13% (-33.34%)

Table 18: SubDefense for I2P-nudity on FMN against UnlearnDiff, with 100 blocked tokens.

Scenario:	FMN	FMN + SubDefense
NoAttack	87.94%	37.59% (-50.35%)
UnlearnDiff	97.87%	45.39% (-52.58%)

Table 19: SubDefense for I2P-nudity on UCE against UnlearnDiff, with 100 blocked tokens.

Scenario:	UCE	UCE + SubDefense
NoAttack	21.98%	13.47% (-8.51%)
UnlearnDiff	78.72%	45.39% (-33.33%)

Table 20: SubDefense for I2P-nudity on SPM against UnlearnDiff, with 100 blocked tokens.

Scenario:	SPM	SPM + SubDefense
NoAttack	55.31 %	34.04% (-21.27%)
UnlearnDiff	91.49 %	58.97% (-32.52%)

# D.3 DEFENDING AGAINST SUBATTACK ON VARIOUS CONCEPTS FOR VARIOUS UNLEARNED MODELS

Apart from the major baseline comparison of defense on UCE, and the defense results against different attacks on ESD presented in the main paper, we provide additional defense results of various concepts and unlearned models against SubAttack in this section. The results are shown in **Tab. 21**, **Tab. 22**, **Tab. 23**, and **Tab. 24** accordingly. Notice that ASR on various concepts is reduced with SubDefense, while ASR reduction on "Van Gogh" is the most significant. It is worth exploring in the future to design new methods and make the defense more robust for other concepts as well.

Table 21: SubDefense for three concepts on ESD against SubAttack, with 100 blocked tokens.

Scenario:	ESD	ESD + SubDefense
Nudity	97.56%	42.33% (-55.23%)
Van Gogh	81%	17% (-64%)
Church	91.33%	40.22% (-51.11%)

#### D.4 DEFENDING RESULTS ON OTHER EXPLORATORY SETTINGS

**Defending against black-box attack.** We also conducted exploratory experiments on defending against Ring-A-Bell (Tsai et al., 2024), a classic black-box attack. Our results in **Tab. 25** show that SubDefense reduces ASR across several unlearned models, including MACE, FMN, SPM, and ESD. These findings suggest that SubDefense can provide robustness in black-box scenarios, although our main focus remains on white-box settings.

Table 22: SubDefense for three concepts on FMN against SubAttack, with 100 blocked tokens.

Scenario:	FMN	FMN + SubDefense
Nudity	100%	62.89% (-37.11%)
Van Gogh	96.33%	22.78% (-73.55%)
Church	82.67%	13.78% (-68.89%)

Table 23: SubDefense for three concepts on UCE against SubAttack, with 100 blocked tokens.

Scenario:	UCE	UCE + SubDefense
Nudity	81.67%	28% (-53.67%)
Van Gogh	93.78%	14.33% (-79.45%)
Church	82.67%	3.22% (-79.45%)

Table 24: SubDefense for three concepts on SPM against SubAttack, with 100 blocked tokens.

Scenario:	SPM	SPM + SubDefense
Nudity	74.89%	50.78% (-24.11%)
Van Gogh	82.78%	12.33% (-70.45%)
Church	84.89%	23.78% (-61.11%)

Table 25: Exploratory defense results against the black-box Ring-A-Bell (Nudity) attack.

Scenarios:	MACE	FMN	SPM	ESD
Ring-A-Bell ASR	11.58%	95.79%	34.74%	57.89%
+ SubDefense	<b>5.26%</b> (k=10)	<b>54.75%</b> (k=100)	<b>14.74</b> % (k=100)	<b>4.21%</b> (k=100)

**Standalone performance of SubDefense.** Although SubDefense was primarily designed as a plugin defense to enhance the robustness of existing unlearned models (similar to RECE operating on UCE), we also explored its effectiveness as a standalone unlearning method. Specifically, we applied SubDefense directly to the original Stable Diffusion (SD) model without any prior unlearning. Results are promising: as shown in **Tab. 26**, SubDefense reduces ASR under both black-box (Ring-A-Bell, Nudity) and white-box (SubAttack, Church) attacks.

Table 26: Standalone performance of SubDefense.

Scenarios:	SD	K=10	K=20	K=50	K=100	K=150	K=200
Ring-A-Bell (Nudity)	97.89%	89.47%	76.84%	60%	38.94%	23.16%	8.42%
SubAttack (Church)	100%	80%	78.33%	55%	46.67%	21.67%	10%

#### E ABLATIONS

#### E.1 ATTACK

**Number of attack tokens.** In practice, we use K=5 to conduct SubAttack as it provides strong attack performance while maintaining computational efficiency. Here, we take ESD as an example to show how ASR varies with K. To conduct ablations more efficiently, we subsample 300 out of 900 prompts for the concepts "church" and "nudity" to study the relationship between ASR and K. Results are presented in **Fig. 12** and **Fig. 13**. The additional attack time per image caused by

each additional token embedding is approximately 10 seconds, which leads to about 3 more hours to attack a single concept having 900 prompts in the dataset. Therefore, considering the needs of attacking multiple concepts and multiple models in practice, we choose K=5 where the ASR is approximately stabilized. For some unique scenarios, users can choose to increase K for higher ASR at a cost of longer computation time.

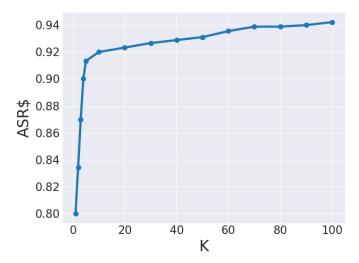


Figure 12: **ASR versus** K when conducting SubAttack on ESD for the concept "church".

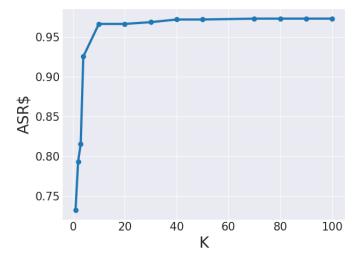


Figure 13: **ASR versus** *K* when conducting SubAttack on ESD for the concept "nudity".

**Orthogonility.** The orthogonality constraint was introduced to encourage diversity among the learned attack embeddings, preventing them from collapsing into a single semantic direction and thereby covering a broader and more effective attack space. To validate this design choice, we conducted an ablation study on the "Nudity" concept. As shown in **Tab. 27**, enforcing orthogonality consistently improves ASR across multiple unlearned models, supporting the effectiveness of this constraint.

**Vocabulary size.** We ablate the vocabulary size used for SubAttack by selecting the top-N CLIP tokens most similar to the target concept. As shown in **Tab. 28**, ASR improves sharply up to 5000 tokens but declines when the vocabulary grows larger. This indicates that 5000 tokens strike the best balance between diversity and optimization feasibility.

Table 27: **Ablation on the orthogonality constraint.** Enforcing orthogonality improves ASR across unlearned models for the "Nudity" concept.

Scenarios:	ESD	FMN	UCE	SPM
With Orthogonality	97.56%	100%	81.67%	74.89%
Without Orthogonality	78.33%	100%	71.67%	63.33%

Table 28: Ablation on vocabulary size. ASR of SubAttack on ESD "Nudity" with different vocabulary sizes.

Vocabulary size	50	500	5000 (default)	10000	Full
ASR	43.33%	81.67%	97.56%	70%	28.33%

#### E.2 Defense

 **Gradual degradation of generation utility with stronger defense.** We show an ablation study on COCO-10k generation CLIP score and FID versus the number of blocked tokens in **Tab. 29** using ESD for the concept of "nudity". We can see that, after the number of blocked tokens surpasses 100, there appears to be a significant harm to the CLIP score and FID. In practice, the number of blocked tokens during defense can be selected to balance good generation quality and low ASR according to one's preference. In this paper, we provide an ablation study on ESD as an example, and report ASR majorly with 20 or 100 blocked tokens for different unlearned models and concepts.

Table 29: SubDefense exhibits gradual degradation of CLIP score and FID when the number of blocked token embeddings increases.

#Blocked Tokens:	0	20	50	100	200	300	350
CLIP Score (†)	30.13	30.02	29.86	29.58	28.54	26.15	24.72
FID (↓)	18.23	19.02	19.09	19.20	20.92	26.42	30.33

More results and discussions on defending against CCE. Defending against CCE is an underexplored problem in the field, where there are no baselines to compare with, to the best of our knowledge. Hence, we show a detailed study on defense against CCE, along with more discussions to support future research. As shown in Tab. 30, different from UlearnDiff, CCE requires a large number of tokens to be blocked if we aim to have low ASR. However, lower ASR achieved by more blocked attack tokens leads to a degradation of generation utility, with an increased FID and a decreased CLIP score, referring to **Tab. 29**. Such a phenomenon indicates that the embedding identified by CCE has a complex association with the target concept, sharing components with a variety of interpretable token embeddings found by our method. This suggests that fully understanding the behavior of CCE requires a deeper analysis of how LDMs interpret and generate concepts other than the current approach we use. For example, currently, the interpretability of retained associations of concepts relies on predefined CLIP vocabularies, which may not capture all implicit or nuanced representations retained in unlearned models. While the above question is beyond the scope of the current work, such insights could inform the development of more robust and versatile defense strategies in the future. With improved understanding of LDMs, future research may come up with more efficient and robust defenses against CCE while preserving model utility.

Table 30: ASR of concept "nudity" on CCE after blocking different numbers of token embeddings.

#Blocked Tokens:	0	100	230	270	320	350	390
CCE ASR	85.11%	75.67%	65.78%	37.44%	28.11%	18.11%	8.89%

## F SPARSITY OF ATTACK TOKEN EMBEDDINGS

Sparsity constraints are widely adopted in prior concept decomposition works - where the linear combination coefficients  $\alpha_i$  are forced to be nearly zeros except for dozens of tokens (usually 20-50). However, in our attacks, where the unlearned diffusion models majorly associate the target concept with a set of implicit tokens, removing such sparsity regularization is helpful, especially for attack token embeddings discovered later in the iterative learning process. Hence, we do not impose a sparsity constraint. Yet, it's interesting to find through our learning that a weaker sparse structure still emerges, and such sparsity gradually decreases as we learn more attack token embeddings through the iterative learning process.

Specifically, for each learned attack token embedding, we normalize  $\alpha = [\alpha_1, \dots, \alpha_N]$  to have a unit norm. Then, we find the index  $i^*$  such that:

$$i^* = \operatorname*{arg\,min}_i i$$
, such that  $\sum_{j=1}^i \alpha_j^2 \ge 0.9$  (5)

Besides, we also count the number of  $\alpha_i$  such that  $\alpha_i \ge 0.01$ . We report the results of the first attack token embedding on ESD for each concept in **Tab. 31**. Notice the size of the CLIP token vocabulary is more than 40000.

Table 31: Sparsity of the learned attack token embeddings.

Concept:	Nudity	Van Gogh	Church
$\overline{i^*}$	1455	668	547
$\#\alpha_i \geq 0.01$	1743	1023	885

During our iterative learning process of a set of tokens for the nudity concept, we observe a decreasing sparsity, as shown in **Tab. 32**. This is intuitive since later attacking requires more complex associations to the target concept.

Table 32: Sparsity of the learned attack token embeddings decreases during the iterative subspace attack process.

#Itrs	1	10	30	50	70	100	130	150	170	200
$\overline{i^*}$	1455									
$\#\alpha_i \ge 0.01$	1743	2019	2078	2009	2206	2298	2328	2368	2358	2326

Furthermore, we visualize the nudity concept attacking results on ESD by selecting only the largest dozens of  $\alpha_i$  within a learned  $\alpha$  and setting other entries as zeros. As shown in **Fig. 14**, we see the nudity concept is gradually enhanced as the number of selected  $\alpha_i$  increases to 1500: the woman generated happens to wear fewer and fewer clothes until she's completely bare.













Figure 14: Attacking the concept nudity on ESD when  $\alpha$  has different numbers of non-zero entries.

## G IMAGE GENERATION QUALITY VISUALIZATION AFTER DEFENSE

In this section, we provide a more detailed study on the generation quality of unlearned models after we plug SubDefense into them. First, we provide more detailed MSCOCO prompts and the generated images of UCE and UCE + SubDefense (with 20 blocked tokens) in Fig. 15, Fig. 16, and Fig. 17. Next, taking UCE and "Van Gogh" as an example, whose attack token embeddings are highly related to "blue" and "star", we study whether SubDefense of "Van Gogh" harms the generation of "blue" and "star" in Fig. 18 and Fig. 19. It turns out that the ability to generate these related concepts is highly preserved, which highlights that subdefense is different from direct token blocking of all related concepts. Instead, SubDefense blocks the composed embeddings, which represent the concept "Van Gogh" more accurately.

Prompt	UCE	UCE + SubDefense
"A man riding a motorcycle with a woman on back of it."		
"An Air Force jet flying in a deep blue sky."		
"A red car sitting on to of a black boat."		

Figure 15: More detailed visualization of COCO generation results with or without SubDefense on the concept nudity.

Prompt	UCE	UCE + SubDefense
"A car that seems to be parked illegally behind a legally parked car."		
"A couple of birds fly through a blue cloudy sky."	* *	<b>*</b> *.
"A room with blue walls and a white sink and door."		

Figure 16: More detailed visualization of COCO generation results with or without SubDefense on the concept Van Gogh.

Prompt	UCE	UCE + SubDefense
"A bike parked next to a cat leaning up against a stone wall."		
"Two giraffes standing next to each other at a zoo."		
"A black and white cat sits in a white sink."		

 $Figure\ 17:\ More\ detailed\ visualization\ of\ COCO\ generation\ results\ with\ or\ without\ SubDefense\ on\ the\ concept\ church.$ 

Prompt	UCE	UCE + SubDefense
"A majestic blue butterfly resting on a flower."		
"A futuristic city glowing with blue neon lights."		
"A peaceful lake reflecting the <b>blue</b> sky."		

 $\begin{tabular}{ll} Figure~18: Visualization~of~"blue"~image~generation~results~before~and~after~defending~"Van~Gogh"~on~UCE. \end{tabular}$ 

Prompt	UCE	UCE + SubDefense
"Bright <b>star</b> in the night sky."		
"Galaxy with many stars."		
"Glowing star-shaped lantern."		

 $\label{thm:condition} \mbox{Figure 19: Visualization of "star" image generation results before and after defending "Van Gogh" on UCE. \\$ 

## H MORE ATTACK VISUALIZATIONS



Figure 20: Visualizing nudity attacking results on ESD.



Figure 21: Visualizing nudity attacking results on FMN.



Figure 22: Visualizing nudity attacking results on UCE.



Figure 23: Visualizing nudity attacking results on SPM.

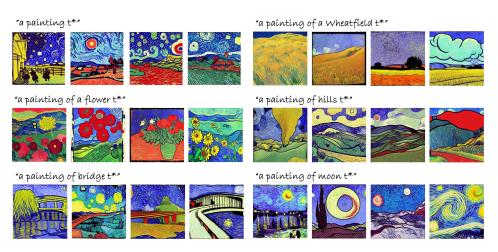


Figure 24: Visualizing Van Gogh attacking results on ESD.



Figure 25: Visualizing Van Gogh attacking results on FMN.



Figure 26: Visualizing Van Gogh attacking results on UCE.

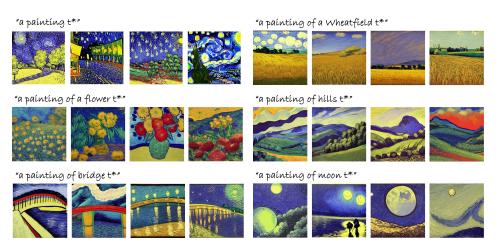


Figure 27: Visualizing Van Gogh attacking results on SPM.



Figure 28: Visualizing church attacking results on ESD.

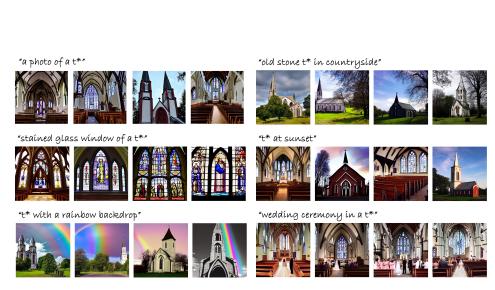


Figure 29: Visualizing church attacking results on FMN.



Figure 30: Visualizing church attacking results on SPM.

## I FUTURE DIRECTIONS

We identify the following future directions. First, future research could explore feature representations in diffusion models beyond the linear structure, which may reveal richer mechanisms underlying unlearning. Second, efficient, adaptive, and automatic methods could be designed to determine not only the number of blocked tokens but also the specific set to block, for example through learned importance scores or attention-based relevance. Third, joint visual—textual embeddings could be investigated to better understand and defend against multimodal jailbreaks. Fourth, as a reference point for defenses against CCE, SubDefense highlights a clear trade-off between robustness and utility; addressing this trade-off remains an important open challenge. Fifth, extending SubDefense beyond CLIP-based architectures is another promising avenue. The core principle of identifying and nullifying harmful semantic directions in the conditional embedding space could be applied to other text encoders or even to models conditioned on alternative modalities. Finally, examining residual associations without relying solely on predefined vocabularies may capture more implicit or nuanced concepts retained in unlearned models, improving interpretability and guiding the development of stronger defenses.