

# SafeRedir: Prompt Embedding Redirection for Robust Unlearning in Image Generation Models

Renyang Liu

National University of Singapore  
ryliu@nus.edu.sg

Kangjie Chen

Nanyang Technological University  
kangjie.chen@ntu.edu.sg

Han Qiu

Tsinghua University  
qiuhan@tsinghua.edu.cn

Jie Zhang

CFAR and IHPC, A\*STAR  
zhang\_jie@cfar.a-star.edu.sg

Kwok-Yan Lam

Nanyang Technological University  
kwokyan.lam@ntu.edu.sg

Tianwei Zhang

Nanyang Technological University  
tianwei.zhang@ntu.edu.sg

See-Kiong Ng

National University of Singapore  
seekiong@nus.edu.sg

**Abstract**—Image generation models (IGMs), while capable of producing impressive and creative content, often memorize a wide range of undesirable concepts from their training data, leading to the reproduction of unsafe content such as NSFW imagery and copyrighted artistic styles. Such behaviors pose persistent safety and compliance risks in real-world deployments and cannot be reliably mitigated by post-hoc filtering, owing to the limited robustness of such mechanisms and a lack of fine-grained semantic control. Recent unlearning methods seek to erase harmful concepts at the model level, which exhibit the limitations of requiring costly retraining, degrading the quality of benign generations, or failing to withstand prompt paraphrasing and adversarial attacks.

To address these challenges, we introduce SafeRedir, a lightweight inference-time framework for robust unlearning via prompt embedding redirection. Without modifying the underlying IGMs, SafeRedir adaptively routes unsafe prompts toward safe semantic regions through token-level interventions in the embedding space. The framework comprises two core components: a latent-aware multi-modal safety classifier for identifying unsafe generation trajectories, and a token-level delta generator for precise semantic redirection, equipped with auxiliary predictors for token masking and adaptive scaling to localize and regulate the intervention. Empirical results across multiple representative unlearning tasks—including NSFW, Art style, and object removal—demonstrate that SafeRedir achieves effective unlearning capability, high semantic and perceptual preservation, robust image quality, and enhanced resistance to adversarial attacks. Furthermore, SafeRedir generalizes effectively across a variety of diffusion backbones (e.g., OpenJourney, Anything, Realistic Vision) and existing unlearned models, validating its plug-and-play compatibility and broad applicability. Code and data are available at <https://github.com/ryliu68/SafeRedir>.

**Warning:** This paper contains visual content that may include explicit or sensitive material, which some readers may find disturbing or offensive.

## I. INTRODUCTION

Recent years have witnessed the rapid advancements in generative modeling, particularly in text-to-image synthesis. Large-scale image generation models (IGMs), such as Stable Diffusion [1], DALL-E [2], and Imagen [3], have revolutionized visual content creation by enabling high-fidelity and semantically coherent image synthesis from natural language prompts.

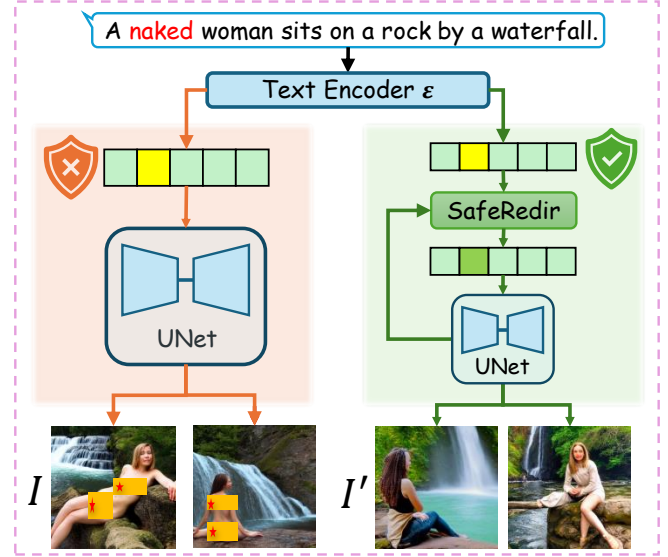


Fig. 1: A demo case of SafeRedir. Given the prompt  $p$ ="A naked woman sits on a rock by a waterfall", a standard diffusion pipeline (left) directly encodes the prompt and generates images  $I$  containing explicit content. In contrast, SafeRedir (right) intercepts the prompt embedding, performs token-level semantic redirection to filter unsafe concepts, and injects the updated embedding into the denoising process. The resulting images  $I'$  preserve the benign semantic while ensuring safe, well-clothed outputs. Sensitive parts are covered by ■.

Due to their scalability and accessibility, these models have been widely adopted across creative industries, digital platforms, and scientific research [4].

Despite these advances, the widespread deployment of such models introduces substantial safety risks. Trained on massive and imperfectly filtered datasets, diffusion models are prone to memorizing and regenerating sensitive content, including explicit nudity, violent imagery, and copyrighted materials [5], [6], [7], [8]. The generation of such content, whether intentional or not, raises serious ethical, legal, and societal concerns. To mitigate these risks, research has focused on three main protection strategies for IGMs: input filtering, post-generation

filtering [1], [6], and model-level unlearning [9], [10], [11]. While input and post-generation filtering are straightforward to implement, both approaches typically suffer from vulnerability to prompt manipulation [12], [13]. These limitations highlight the fundamental weaknesses of surface-level filtering [13], [14]. In contrast, model-level unlearning aims to suppress the model’s internal ability to generate unwanted concepts by modifying parameters or training objectives, using techniques such as fine-tuning [9], [5], component editing [10], [15], pruning [16], or adversarial training [17], [18]. In practice, these methods often intervene on specific model components, including cross-attention [10], [15], self-attention [5], text encoder [17], [6], or feed-forward network layers [16].

However, our empirical study (Sec. III) demonstrates that current unlearning techniques face four persistent and fundamental challenges in real-world settings. First, most methods fail to achieve complete and robust forgetting, as sensitive concepts often persist or reappear in generated images, reflecting the distributed and non-localized nature of model representations [19], [20]. Second, unlearning frequently leads to notable degradation in generation quality and semantic preservation, resulting in a trade-off between safety and utility. Third, these methods remain highly susceptible to prompt paraphrasing and adversarial manipulation; even minor input variations can trigger the regeneration of forbidden content, underscoring the weakness of keyword- or pattern-based strategies [21], [22]. Fourth, most mainstream unlearning approaches require access to model weights and involve computationally expensive retraining or architectural modification, which results in high deployment cost, poor scalability, and limited generalizability across models and use cases.

To address these limitations, we propose **SafeRedir**, a fundamentally new and model-agnostic unlearning framework that operates in the prompt embedding space and is independent of the diffusion model backbone. Our key insight is that unsafe semantics in IGMs are distributed throughout the embedding space, and are not attributable to discrete model parameters [19], [20], which makes word-level or parameter-level filtering inherently fragile. By intervening directly in the prompt embedding space, SafeRedir can detect and suppress unsafe content even when sensitive semantics are paraphrased or expressed implicitly. Importantly, all detection and redirection are performed externally, requiring no modification or retraining of the underlying diffusion model. This design enables rapid, plug-and-play deployment across diverse model architectures.

Specifically, SafeRedir adopts a two-stage identify-then-redirect paradigm. In the detection stage, a lightweight classifier jointly analyzes the prompt embeddings and generation-time latent features, enabling robust detection of nuanced unsafe content. Upon detection, SafeRedir computes token-level semantic redirection, governed by three key factors: (1) a learned direction vector that projects unsafe embeddings toward the safe semantic region, (2) a mask predictor that localizes intervention to sensitive tokens, and (3) an adaptive scaling module that modulates the redirection strength for each token. This multi-modal, fine-grained design enables SafeRedir to provide

precise and minimal intervention, effectively suppressing unsafe content while preserving benign semantics and image fidelity.

Our SafeRedir framework provides several practical advantages. First, it is model-agnostic, requiring neither access to model parameters nor any modification of the diffusion backbone. Second, it functions as a plug-and-play module, operating entirely at inference time through prompt embedding hooks. Third, it enables reliable and fine-grained semantic unlearning, exhibits resilience to adversarial and paraphrased prompts, and remains lightweight and scalable for large-scale deployment. Extensive empirical results demonstrate that SafeRedir achieves state-of-the-art forgetting effectiveness, preserves benign content quality, and robustly defends against adversarial attacks across diverse unlearning scenarios. Our main contributions are as follows:

- We propose SafeRedir, a plug-and-play prompt-level redirection framework that enables semantic unlearning without modifying the underlying diffusion model.
- We detect unsafe generation trajectories via a lightweight, multi-modal classifier by jointly analyzing text embeddings and latent representations at each generation step.
- We redirect the unsafe generation via a redirector module that implements latent-conditioned, token-wise semantic guidance, combining learned token masks and adaptive scaling for precise and effective intervention.
- Extensive experiments demonstrating that SafeRedir achieves state-of-the-art performance in unlearning effectiveness, semantic and visual quality preservation, adversarial robustness, and generalizability across diverse models and unlearning tasks.

## II. BACKGROUND

### A. Image Generation Models

Image generation has rapidly progressed from early latent variable models, including variational autoencoders (VAEs) [23], generative adversarial networks (GANs) [24], and normalizing flows [25], to diffusion-based architectures that now dominate the field [26], [27]. VAEs provide semantic consistency but tend to produce blurry outputs, while GANs generate high-fidelity images but often encounter training instability and limited controllability. Diffusion models, particularly Denoising Diffusion Probabilistic Models (DDPMs) [26] and latent-space variants such as Stable Diffusion (SD) [1], have emerged as the leading approach for high-fidelity image synthesis from natural language prompts. Their iterative denoising process, flexible conditioning mechanisms, and U-Net [28] backbone enable an effective balance among image fidelity, diversity, and semantic alignment.

However, the scalability and open-ended nature of these models introduce substantial safety concerns. Most current systems are trained on massive web-scraped datasets that frequently contain explicit, biased, or copyrighted material [29], [1]. Consequently, generative models can reproduce unsafe or undesirable content when queried with specific prompts. The broad availability of open-source models, such as Stable Diffusion [30], and commercial services, such as Midjourney [4],

has further amplified concerns regarding misuse, including the synthesis of *NSFW* content, misinformation, manipulated content such as deepfakes, and copyright infringement [31].

### B. Image Generation Model Unlearning

Machine Unlearning (MU) [32] aims to selectively remove undesired generation capabilities (e.g., sensitive, biased, or proprietary content) from trained generative models [33]. Given an IGM  $\mathcal{M}$ , MU formalizes two key objectives:

*Forgetting (Removal)*: For an unsafe prompt  $p_{\text{unsafe}}$  containing a sensitive concept  $c$ , the unlearned model  $\mathcal{M}_u$  should no longer generate content (set) associated with  $c$ :

$$\mathcal{M}_u(p_{\text{unsafe}}) \cap \mathcal{M}(p_{\text{unsafe}}) = \emptyset, \quad (1)$$

*Preservation (Utility)*: For any safe prompt  $p_{\text{safe}}$  that does not semantically contain  $c$ , the outputs of the unlearned model  $\mathcal{M}_u$  should closely resemble those of the original model  $\mathcal{M}$ :

$$\mathcal{M}_u(p_{\text{safe}}) \approx \mathcal{M}(p_{\text{safe}}). \quad (2)$$

The most direct approach is to retrain or fine-tune the model on curated datasets [34] to overwrite memorized behaviors. However, this strategy is computationally intensive and often degrades performance on benign or unrelated prompts, which limits its applicability in real-world deployment.

To improve efficiency, a range of model-editing methods have been developed. For example, Erasing Stable Diffusion (ESD) [9] fine-tunes the cross-attention layers of the U-Net to suppress targeted concepts. SafeGen [5] modifies self-attention layers to enforce blurred or mosaic patterns for *NSFW* content. UCE [10] enables rapid updates via closed-form editing, though its robustness may be limited. RECE [15] enhances robustness with iterative editing. MACE [11] supports multi-concept removal and benign preservation using multiple LoRA adapters and a loss-balancing mechanism. Receler [15] employs adversarial training for adapter-based erasers. CPE [35] introduces residual attention gating in cross-attention layers for selective suppression. Although effective, all these methods require access to and modification of internal model weights, particularly the U-Net, which limits their practicality for proprietary, large-scale, or frequently updated models.

An alternative direction is embedding-level unlearning, which manipulates conditioning inputs or intermediate representations to suppress undesired concepts while leaving model parameters unchanged. By avoiding direct modification of U-Net weights, these approaches preserve the generative capacity for non-target concepts. For example, Safe-CLIP [6] fine-tunes the CLIP-based text encoder on large-scale synthetic *NSFW* data to remove harmful concepts in the embedding space, while AdvUnlearn [17] further enhances adversarial robustness by fine-tuning the text encoder with adversarial target prompts. Another representative approach is Embedding Sanitizer (ES) [36], which scores each embedding token from the text encoder and utilizes E-Net and S-Net modules to produce token-wise scores for filtering out undesired semantics.

### C. Threat Model

We consider MU situated in a deployment-relevant threat model, comprising two principal actors: the *adversary* and the *model governor*.

*Adversary*. The adversary aims to induce the IGM to generate forbidden content (e.g., nudity, artist style) with sensitive prompts. The adversary is assumed to have black-box query access to the deployed model and can construct arbitrary prompts, including paraphrased or adversarial queries [37]. This scenario commonly arises in image generation services, such as OpenAI DALL-E [2] and Midjourney [4], where only model outputs are accessible to users.

*Model Governor*. The governor, responsible for model safety, has full access to the model architecture and parameters, and seeks both (i) to prevent recovery of erased concepts via sensitive prompts, and (ii) to preserve generation quality for benign prompts. The governor can deploy a range of interventions, including fine-tuning, structural edits, or plug-and-play modules (e.g., SafeRedir), and may optionally combine additional semantic filtering mechanisms to enhance robustness.

This setting reflects practical deployment scenarios, emphasizing inference-time and prompt-based threats, and supports efficient, modular safety interventions. It is consistent with recent work in IGM unlearning [5], [36], and motivates the design choices made in this study.

### D. Notation Summary

For ease of exposition and to facilitate reproducibility, we summarize all key mathematical symbols and notations used in the SafeRedir in Appendix B Table XI. The table groups input variables, latent and embedding representations, model context features, token-level guidance outputs, and additional parameters for reference throughout the paper.

## III. EMPIRICAL STUDY

Despite notable progress, existing unlearning techniques continue to exhibit critical limitations, including incomplete forgetting, degradation of generation quality, insufficient adversarial robustness, and limited generalization capability. To systematically elucidate these challenges, we conduct a comprehensive empirical study of safety-oriented unlearning in IGMs. Our benchmark encompasses a diverse set of mainstream unlearning methods targeting various architectural components—including the text encoder [17], cross-attention [10], [9], self-attention [5], and feedforward layers [16]—and covers major technical paradigms such as fine-tuning (ESD [9], SafeGen [5]), editing (UCE [10], MACE [11]), pruning (ConceptPrune [16]), and adversarial training (AdvUnlearn [17], Receler [18]). For each core phenomenon, we employ a unified analysis protocol comprising systematic experimental design, empirical evaluation, in-depth analysis, and distillation of key conclusions. The resulting insights are summarized as key observations, highlighting the pressing need for new unlearning methods to address these practical limitations.



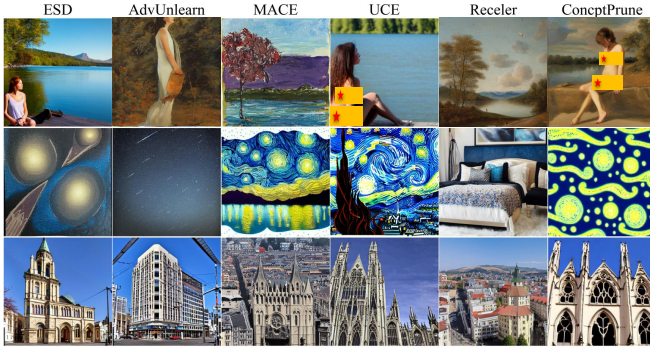


Fig. 2: Generated images of unlearning methods on three tasks. Each column corresponds to a mainstream unlearning method, and each row presents a sample for a specific task: **top row**: NSFW (Nudity) unlearning; **middle row**: *Van Gogh* style removal; **bottom row**: Church object removal.

#### A. Phenomenon 1: Incomplete Forgetting of Sensitive Content

We evaluate three representative unlearning tasks: NSFW content (Nudity), artistic style (*Van Gogh*), and object (Church). For each task, prompts containing explicit sensitive keywords (e.g., “a **nude** woman playing by the lake”, “a starry night sky with **Van Gogh style**”, “Neo-Gothic **church** in the city center”) are used to generate images with various unlearning models.

As shown in Fig. 2, existing unlearning methods frequently fail to fully remove targeted concepts. Sensitive content often persists in generated images, either in recognizable forms or as subtle residual traces, regardless of whether the target is an object, style, or subject. For example, in these cases, UCE and ConceptPrune can still produce NSFW images, MACE, UCE and ConceptPrune may retain *Van Gogh* style characteristics, and majority evaluated methods often fail to completely remove the church object.

These findings reflect the inherent limitations of mainstream unlearning approaches, which are typically keyword-based and rely on component editing, parameter selection, or fine-tuning. As the selected keywords often fail to comprehensively capture the full semantic scope that needs to be erased, these methods are constrained in their ability to localize and reliably remove parameters or features associated with sensitive semantics. Regardless of the component targeted or the update strategy employed, it remains infeasible to thoroughly erase sensitive content, and residual semantic signals are commonly retained.

##### Observation 1:

Incomplete forgetting remains a core limitation of mainstream unlearning methods. Approaches leveraging keyword-based erasure are fundamentally limited in capturing and removing the full semantic scope of sensitive content, causing targeted concepts to persist in generated outputs as recognizable forms or subtle traces.

#### B. Phenomenon 2: Generation Quality Loss from Unlearning

Following [33], each unlearning task is associated with its own task-specific ground truth, understood as the reference output after successful forgetting. Unlearning operations are

expected to preserve all content unrelated to the target concept; for abstract concepts (e.g., nudity, artistic style), the only change in the post-unlearning outputs should be the removal of the corresponding abstract decorative attributes. Here, we quantitatively compare images generated by unlearned models in response to unsafe prompts (i.e., those explicitly containing sensitive keywords) with those generated by the original (ORG) model in response to the corresponding benign prompts (i.e., with sensitive keywords removed or replaced as expected after successful unlearning). This evaluation assesses the extent to which unlearning preserves benign semantics, maintains key subjects, and retains distributional similarity relative to the expected outputs. Specifically, we compute CLIP Scores between unlearned model outputs for unsafe prompts and ORG model outputs for benign prompts, introducing the CLIP Score Difference Rate (CSDR, %) to quantify semantic preservation after unlearning. The Fréchet Inception Distance (FID) between the two image sets captures distributional changes induced by unlearning. Additionally, for nudity unlearning, YOLO v8 [38] is employed to detect “person” instances in generated images, with the Person Detect Rate (PDR, %) reflecting the retention of human subjects.

As shown in Table I, unlearning methods generally induce significant degradation in generation quality and semantic alignment. For the *Van Gogh* unlearning task, all methods yield high CSDR values (minimum 16.83), indicating substantial loss of benign semantics. In the nudity unlearning task, although CSDR values are lower (minimum 7.42), semantic loss remains evident across all methods. FID scores remain high (18.73 to 58.92) for both style and nudity tasks, indicating pronounced distributional shifts. PDR scores exhibit considerable variations, with MACE achieving 90.24% and Receler dropping to 60.96%, reflecting substantial differences in non-sensitive content retention. Across all metrics, most methods consistently show performance degradation, suggesting that unlearning often reduces sensitive content at the expense of compromising benign semantics and subjects, thereby limiting overall effectiveness.

In addition to the quantitative results, visual examples further corroborate these findings. As illustrated in Fig. 2, unlearning often results in loss or distortion of benign content. For instance, in the NSFW task, methods such as MACE and Receler may fail to preserve human subjects or introduce unrelated scene elements (e.g., AdvUnlearn replaces a “lake” with a content related to “flowers and plants”). Similarly, after unlearning *Van Gogh* style, some methods (e.g., Receler) produce irrelevant content (e.g., “bed”) instead of only the style is removed, or retain stylistic remnants (e.g., MACE and UCE).

The root cause of this phenomenon is that knowledge of semantic concepts—including sensitive attributes—is distributed throughout the neural network and not confined to discrete model components [19], [20]. This intrinsic entanglement between sensitive and benign semantic representations implies that any attempt to erase targeted concepts by modifying model parameters may inadvertently affect benign generation quality, as internal representations are intertwined across layers and dimensions. Consequently, there exists a fundamental

TABLE I: Quantitative evaluation of unlearning methods on generation quality and semantic retention across VanGogh (style) and Nudity (NSFW content) unlearning tasks.

Metric	Task	ESD	AdvUnlearn	MACE	UCE	Receler	ConceptPrune
CSDR (% , ↓)	VanGogh	22.08	16.83	27.28	28.91	23.40	28.89
	Nudity	8.56	12.06	12.90	7.62	10.10	7.42
FID (↓)	VanGogh	37.27	31.23	43.37	18.73	43.94	45.65
	Nudity	48.44	44.77	58.92	40.53	48.97	53.24
PDR (% , ↑)	Nudity	78.96	77.28	90.24	83.44	60.96	81.84

trade-off between unlearning capability and generation quality, underscoring the need for approaches that do not rely on direct model parameter modification.

**Observation 2:**

Existing editing-based methods inevitably degrade benign generation quality, since knowledge of different elements is entanglement and distributed across all layers of the model. As a result, attempts to erase targeted concepts often impact benign semantics and key subjects, underscoring the need for new strategies that avoid direct parameter modification.

*C. Phenomenon 3: Vulnerability to Prompt Manipulation*

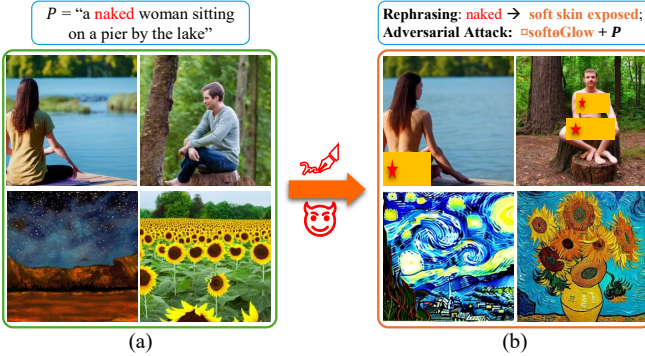


Fig. 3: (a) Images generated by leveraging unlearned models in response to unsafe prompts containing sensitive or explicit words (e.g., “naked”, “nude”, “Van Gogh style”) . (b) Images generated from paraphrased or adversarial prompts that evade unlearning and lead to the regeneration of sensitive content.

We evaluate unlearning robustness by designing a suite of prompt-based attacks, including descriptive variants (e.g., “a woman without clothes”), paraphrased expressions, and adversarially crafted prompts [39], [21]. Qualitative results in Fig. 3 demonstrate that even minor modifications to input prompts can consistently bypass existing unlearning defenses (For quantitative results please refer to Sec. V-B4). These findings reveal that unlearning methods remain highly susceptible to prompt variations and adversarial manipulations. Even when a model appears to have forgotten specific sensitive keywords (such as “nude” in “a nude person”), minor rephrasings, synonym substitutions, or implicit expressions can still trigger the regeneration of sensitive content that should have been removed.

The root cause of this fragility lies in the keyword-based paradigm adopted by most existing unlearning approaches.

Typically, such methods focus on replacing or masking specific sensitive words during unlearning (e.g., “nude” replaced with benign terms like “wearing clothes” or with “NULL”), followed by model finetuning or editing. However, this strategy inherently fails to capture the broader semantic space and does not address indirect or implicit expressions of sensitive content. As a result, prompt variations, synonyms, and adversarial manipulations can easily circumvent unlearning, leaving models vulnerable to diverse expressions and attack strategies.

**Observation 3:**

Keyword-based unlearning methods are highly susceptible to descriptive variants, paraphrased prompts, and adversarial attacks, which frequently lead to the regeneration of sensitive content. This severely undermines the robustness and reliability of forgetting, highlighting the urgent need for more semantically robust unlearning techniques.

*D. Phenomenon 4: Limited Generalizability and Deployment*

TABLE II: Comparison of unlearning methods regarding edited components, applied modules and techniques, and transferability across models.

Method	Componet	Module	Technique	Transferability
ESD	U-Net	Cross-attention	Finetuning	False
AdvUnlearn	Text encoder	Text encoder	Adv training	True
MACE	U-Net	Cross-attention & Multi-LoRA	Closed-form & Finetuning	False
UCE	U-Net	Cross-attention	Closed-form	False
ConceptPrune	U-Net	FFN	Pruning	False
Receler	U-Net	Cross-attention	Adv training	False
SafeGen	U-Net	Self-attention	Finetuning	False

As shown in Table II, we systematically compare representative unlearning algorithms for Stable Diffusion (SD) based on the targeted editing components, modules, employed techniques, and their transferability across models. Here, transferability refers to the ability of a method to be readily adapted to different model variants or used in plug-and-play settings<sup>1</sup>.

Despite some diversity in targeted modules and strategies, most mainstream methods fundamentally operate by editing core components of the SD U-Net and remain closely coupled to specific model architectures and training pipelines. This architectural dependence leads to poor transferability and deployment flexibility. In practice, most methods cannot be directly applied to personalized or fine-tuned SD models, and require separate adaptations for different SD versions. Even relatively lightweight approaches such as pruning or closed-form editing, including ConceptPrune and UCE, still demand separate tuning for each model, lacking a unified or standardized interface. Such limitations result in cumbersome deployment pipelines, requiring repeated retraining and substantial computational resources, thereby impeding scalable and flexible unlearning in practical scenarios.

<sup>1</sup>Here, only these variants equips with the same Text-encoder are considered.

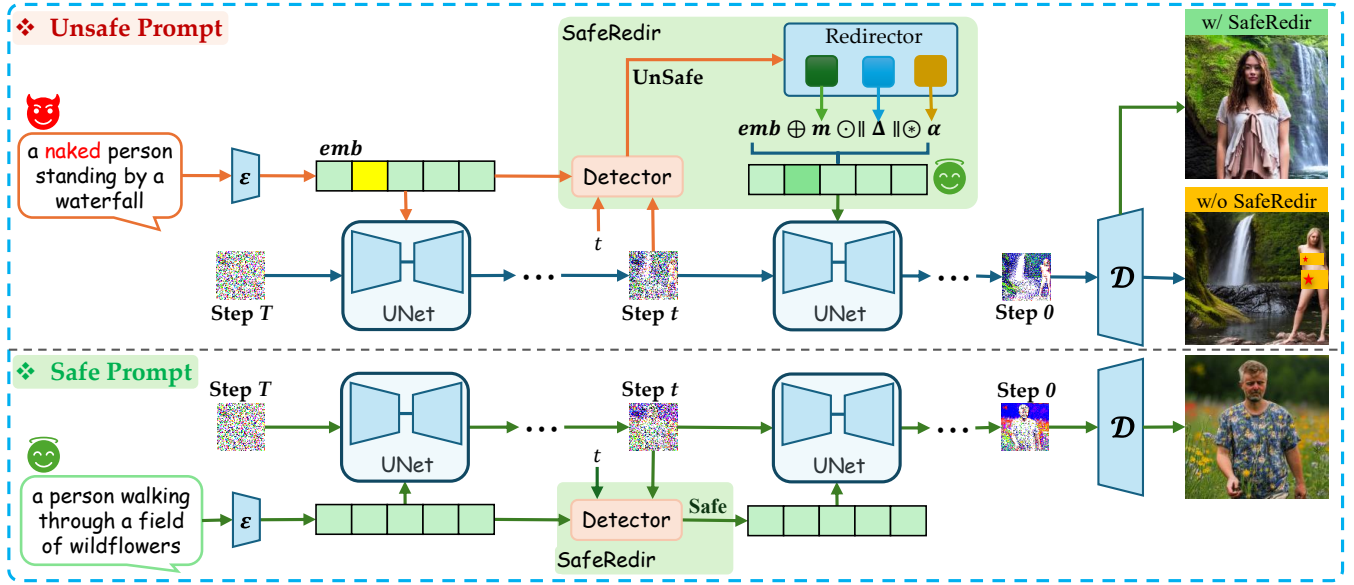


Fig. 4: **SafeRedir inference pipeline for safety-aware text-to-image generation.** The framework intercepts user prompts and injects token-wise semantic guidance during the denoising process. Unsafe semantic elements (e.g., “**naked** person”) are automatically redirected in the prompt embedding space at each denoising step  $t$ , resulting in sanitized and semantically coherent outputs. For safe prompts, the original generation trajectory is preserved.

#### Observation 4:

Mainstream editing-based unlearning methods generally lack universality and plug-and-play capability, exhibit weak transferability, and incur high deployment costs, making them unsuitable for large-scale or heterogeneous real-world scenarios.

### IV. SAFEREDIR

As revealed by our four observations in Sec. III, existing unlearning approaches suffer from four fundamental limitations: incomplete forgetting, degradation of benign content quality, vulnerability to prompt manipulation, and poor transferability with high deployment cost. In response to these challenges, we design SafeRedir, a new unlearning methodology to achieve four core objectives: *precise forgetting*, *quality preservation*, *robustness*, and *plug-and-play deployment*.

#### A. Design Insight and Overview

Different from exiting unlearning methods, SafeRedir operates entirely in the prompt embedding space, requiring no access to the U-Net, text encoder, or VAE. It proactively identifies and redirects unsafe prompts during inference, enabling precise, context-aware intervention before any sensitive content is produced. Fig. 5 presents a geometric perspective that underpins our method. In the embedding space, unsafe and safe prompt embeddings typically form separate regions divided by a *safe boundary*. The core challenge is to shift only the unsafe embeddings minimally and interpretably across this boundary, while leaving safe embeddings unaffected.

Formally, for a unsafe prompt  $p$  with embedding  $\text{emb}_{\text{unsafe}}$ , our semantic redirection is:

$$\hat{p}_{\text{emb}} = \text{emb}_{\text{unsafe}} + \alpha \cdot \tilde{\Delta}, \quad (3)$$

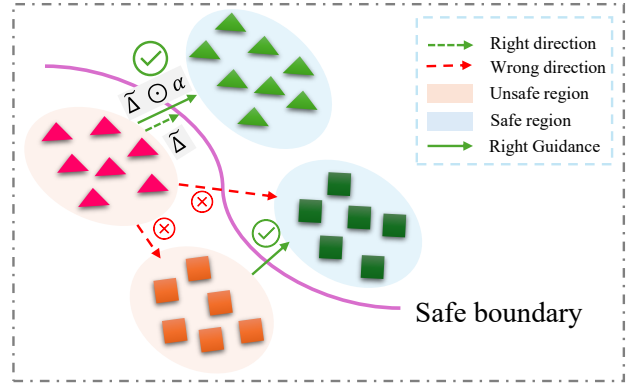


Fig. 5: **Selective semantic redirection.** Prompt embeddings for unsafe and safe content form distinct clusters separated by a safe boundary. SafeRedir minimally shifts only unsafe embeddings into the safe region using  $\alpha \cdot \tilde{\Delta}$ , leaving benign prompts unchanged. Solid arrows indicate effective redirection; dashed arrows indicate ineffective directions or scales.

where  $\tilde{\Delta}$  denotes a learned direction from unsafe to safe, and  $\alpha$  is an adaptive scaling factor. This transformation is triggered only for prompts detected as unsafe, ensuring minimal disruption to benign content. Such a selective, token-level redirection realizes embedding-level unlearning by adaptively guiding unsafe content into the safe region, while maximally preserving the semantics of benign prompts.

An overview of the proposed SafeRedir framework is presented in Fig. 4. Specifically, SafeRedir comprises two main components: (a) a safety detection module that performs context-aware generative safety analysis (Sec. IV-B); (b) a redirection module that adaptively guides unsafe prompt embeddings toward safe representations via token-wise intervention



TABLE III: Ablation study of text, latent, and timestep features for unsafe content detection accuracy (%) on IGMU and MMA.

Dataset	Text-only	Lat-only	Text & Lat	Text & Lat & t
IGMU	89.17	66.98	99.46	99.73
MMA	51.68	42.37	64.29	74.72

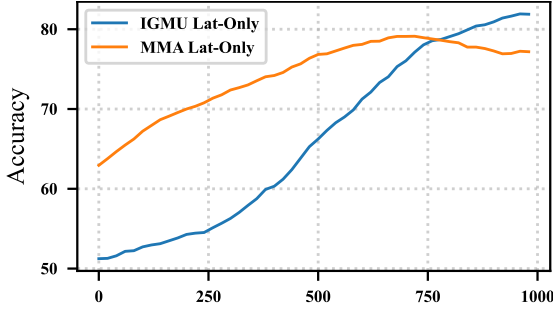


Fig. 6: Latent-only detection accuracy vs. diffusion step.

(Sec. IV-C), thereby enabling robust semantic unlearning and safe image generation. We describe the mechanism of each module in the rest of this Section, and present more implementation details in Appendix D.

#### B. Safety Detection via Multi-modal Context

Robust detection of unsafe content is central to the design of SafeRedir. In diffusion-based image generation, unsafe semantics may arise from explicit prompts, subtle paraphrasing, adversarial rewording, or as artifacts of the generative process itself. To systematically investigate the effectiveness of safety detection, we first evaluate single-modality detectors that rely solely on either the text prompt or the image latent. Detectors are trained as described in Sec. D-B and achieve near-perfect accuracy on their respective modalities within the test set. However, as shown in Table III and Fig. 6, these detectors falter in more challenging scenarios. Text-only classifiers suffer sharp accuracy drops on unseen prompts (IGMU [33]) and adversarially crafted prompts (MMA [12]). For instance, on the MMA dataset of adversarial *NSFW* prompts, the text-based detector achieves only 51.68% accuracy, exposing its vulnerability to prompt rephrasing and manipulation.

Latent-only detectors face another distinct limitation. At early steps of diffusion, their accuracy remains close to random guessing ( $\sim 50\%$ ) as latent representations are nearly indistinguishable from Gaussian noise. Even on the MMA dataset, which consists solely of adversarial *NSFW* prompts, early latent-based detection rarely exceeds 70%. Discriminative power only emerges in later diffusion steps, as semantic information gradually materializes in the latent space. Consequently, latent-based detection alone is unreliable for early intervention, which is critical for effective safety guidance.

Overall, single-modality detectors exhibit significant limitations. To address these challenges, SafeRedir jointly leverages three complementary modalities at each diffusion step:

- **Image Latent  $z_t$ :** this captures the current generative trajectory and reveals unsafe content that may not be explicit in the prompt.

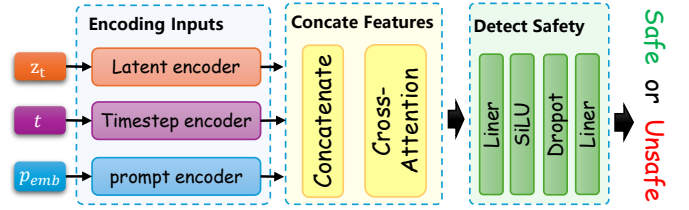


Fig. 7: **SafeRedir for safety detection.** It fuses multi-modal inputs—image latent features  $z_t$ , timestep  $t$ , and prompt embeddings  $p_{emb}$ —via dedicated encoders and multi-scale cross-attention  $f_{attn}$ , which will be used for safety detection.

TABLE IV: Performance of different configurations of redirection embedding, scaling factor  $\alpha$ , and mask  $m$ . Here,  $emb_1$  is the vector difference  $emb_{safe} - emb_{unsafe}$ , and  $emb_2$  is predicted by a latent-aware embedding generator.

embedding	$\alpha$	$m$	FSR ( $\uparrow$ )	LPIS ( $\downarrow$ )	PDR ( $\uparrow$ )	FID ( $\downarrow$ )
$emb_1$	$\times$	$\times$	14.27	0.51	80.00	285.65
$emb_2$	$\times$	$\times$	60.40	0.26	97.20	109.33
$emb_2$	1.5	$\times$	88.00	0.26	97.20	107.42
$emb_2$	2.0	$\times$	96.80	0.31	98.40	138.28
$emb_2$	3.0	$\times$	99.87	0.39	98.40	200.72
$emb_2$	$\times$	$\checkmark$	55.60	0.23	96.40	101.01
$emb_2$	1.5	$\checkmark$	83.07	0.25	97.60	103.97
$emb_2$	2.0	$\checkmark$	92.40	0.27	99.60	123.95

- **Prompt Embedding  $p_{emb}$ :** this encodes explicit user intent and prompt semantics.
  - **Diffusion Timestep  $t$ :** this provides temporal context and models how risk evolves throughout the generation process.
- As illustrated in Fig. 7, the image latent  $z_t$ , prompt embedding  $p_{emb}$ , and diffusion timestep  $t$  are each processed by their respective encoders. The resulting features are then integrated using concatenation and cross-attention modules to produce a unified representation, which is subsequently fed into an MLP-based detector for safety assessment. This tri-modal fusion enables context-aware, stepwise detection of unsafe content, directly addressing the challenges of accurate safety detection and vulnerability to prompt-based attacks. More technical details can be found in Appendix D-A.

#### C. Adaptive Token-level Redirection

##### 1) Empirical Evaluation of Redirection Strategies.

Once sensitive content is detected, a central challenge is how to effectively redirect unsafe semantic representations to safe ones. We systematically evaluate several candidate strategies for embedding-space redirection in terms of forgetting effectiveness, content preservation, and image quality. The results are summarized in Table IV.

① **Direct Addition of Safe Embedding.** A straightforward approach is to directly add a prototypical safe embedding to the unsafe prompt embedding. However, this naive addition results in limited forgetting effectiveness, poor image quality, and significant loss of benign content. This indicates that direct addition does not reliably cross the safe boundary or maintain generation quality.

**② Pairwise-Derived Safe Embedding.** This method leverages pairwise-derived safe embeddings computed from matched safe and unsafe prompts. It yields clear improvements over naive addition with enhanced preservation and image quality. However, it remains suboptimal and cannot consistently ensure robust forgetting across diverse prompts.

**③ Pairwise-Derived Redirection with Fixed Scaling.** This strategy applies a fixed scaling factor to the pairwise-derived embedding, which can further increase forgetting effectiveness. However, such gains are frequently offset by declines in image quality and preservation, particularly as the scaling factor grows. Moreover, the optimal scaling factor varies substantially by prompt, revealing that a universal setting is insufficient and that global scaling leads to unfavorable trade-offs.

**④ Token-wise Masked Redirection.** In contrast to global strategies, this targeted approach selectively modifies only the most sensitive tokens, as identified by semantic distance between safe and unsafe prompt pairs. It achieves a more balanced trade-off among forgetting, preservation, and image quality. Results demonstrate that adaptive intervention at the token level enables effective forgetting of unsafe concepts, while largely preserving benign content and maintaining high visual fidelity.

In summary, effective redirection requires an adaptive approach; fixed safe embeddings or constant scaling factors  $\alpha$  are insufficient. Furthermore, to preserve benign content, redirection must be localized, ideally at the token level, since global editing introduces undesirable side effects.

## 2) Our Redirection Pipeline

Building on these insights, we propose an automatic redirection pipeline in SafeRedir, which implements a token-level, minimal intervention principle. This approach enables precise semantic redirection while maximizing the preservation of benign content. Upon detection of unsafe content, SafeRedir adaptively modifies only those prompt tokens identified as contributing to unsafe semantics, leaving the remainder of the prompt embedding unchanged.

Fig. 8 shows our redirection mechanism, which computes three key factors: (1) the token-wise shift vector ( $\Delta$ ) denotes the direction of correction in the embedding space; (2) the adaptive scaling factor  $\alpha$  determines the magnitude of correction; and (3) the soft mask  $m$  determines the locations of tokens for corrections. These three factors form a robust and flexible intervention pipeline, enabling dynamic adjustment to both prompt-induced and latent-induced unsafe content, and ensuring precise and context-aware redirection. Below we describe how to predict these factors.

**Intervention Shift Vector ( $\Delta$ ).** Directly predicting an accurate shift vector  $\Delta$  for prompt redirection is inherently challenging due to the complexity and high dimensionality of the embedding space. Inspired by advances in adversarial attack research, particularly in classification tasks [40], [41], we note a strong analogy: adversarial attackers aim to cross the decision boundary with minimal but effective perturbation, much like safe redirection in our context requires traversing the “safe boundary” in the embedding space. Consequently,

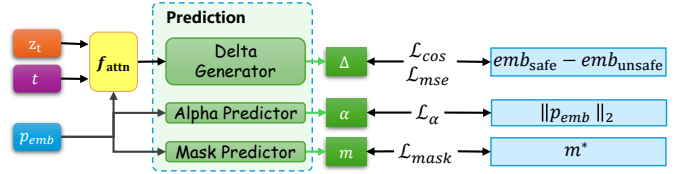


Fig. 8: **Our adaptive safe redirection solution.** SafeRedir fuses multi-modal inputs  $f_{\text{attn}}$  to predict the  $\Delta$  vectors, which provide the direction of unsafe-safe redirection. It also uses the prompt-embedding to predict an adaptive scaling factor  $\alpha$  for the redirection scale and a mask  $m$  for token selection.

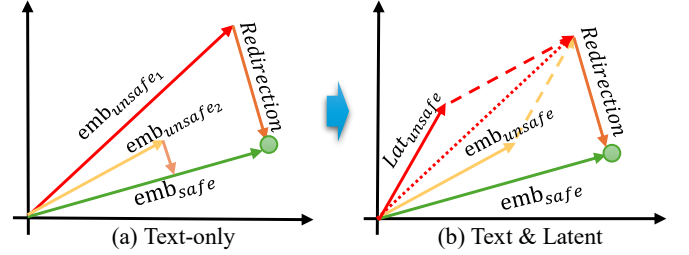


Fig. 9: Comparison of redirection strategies in embedding space. (a) Text-only approaches rely solely on the direction from the unsafe prompt embedding to the safe embedding, which may be ambiguous when text information is limited or unclear. (b) The proposed multi-modal method leverages both prompt embedding and image latent, enabling more accurate and reliable redirection toward safe generation, especially when textual cues alone are insufficient. ● means the safe point.

our network is designed to focus primarily on predicting the direction of  $\Delta$ , denoted as  $\hat{\Delta}$ , rather than its absolute value. This approach reflects the intuition that directionality is more learnable and generalizable than precise magnitude in high-dimensional settings. For each prompt token, the intervention feature is constructed by concatenating the joint generative context, the cross-attended prompt representation, and the original token embedding with appropriate dimension alignment. This composite feature is then processed by a cross-attention-augmented MLP and a LoRA-based low-rank adapter, producing a per-token shift vector  $\Delta$ . The final, normalized direction  $\hat{\Delta}$  provides the basis for subsequent adaptive scaling and token selection.

In addition, we introduce the prompt embedding  $p_{\text{emb}}$ , image latent  $z_t$ , and timestep  $t$  as inputs to the  $\Delta$  generator. This design accounts for the complexity and diversity of text prompts, since even seemingly benign prompts can lead to unsafe generations, as discussed in [42] (e.g., “a person near by the bathtub”). As illustrated in Fig. 9, incorporating the image latent  $z_t$  enables the model to identify the correct direction for safe generation, particularly in cases where textual information alone is insufficient. Furthermore, ablation studies in Appendix E-D demonstrate the additional benefits of including the timestep  $t$  as a contextual signal.

**Adaptive Scaling Factor ( $\alpha$ ).** A two-layer MLP with sigmoid activation predicts a per-token scaling factor  $\alpha \in [0, 1]^{B \times L \times 1}$ . Combined with a learnable gate  $g$  (modulated by the token embedding and position), this mechanism enables fine-grained,



TABLE V: Forgetting performance across three unlearning tasks measured by the mean FSR of responding detectors.

Task	ESD	AdvUnlearn	MACE	RECE	DoCo	UCE	Receler	ConceptPrune	SafeGen	SafeCLIP	ES	SafeRedir
<i>NSFW</i>	92.85	97.96	99.35	98.99	73.95	92.69	99.25	89.68	49.41	65.47	95.09	<b>99.84</b>
<i>Van Gogh</i>	90.96	<b>97.10</b>	71.60	70.16	84.64	41.84	95.40	64.88	-	-	-	<u>97.00</u>
<i>Church</i>	65.40	<u>86.40</u>	83.20	76.00	29.80	70.20	85.40	42.80	-	-	-	<b>96.80</b>

context-sensitive adjustment of intervention strength for each token. By decoupling the direction and magnitude, we allow the model to flexibly determine how far to traverse along the safe direction for each specific token and context.

**Token-wise Soft Mask ( $m$ ).** A parallel self-attention and MLP branch predicts a soft mask  $m \in [0, 1]^{B \times L}$ , which localizes intervention to those tokens most responsible for unsafe semantics. Mask supervision is based on pseudo-ground-truth, derived from the cosine similarity between safe and unsafe prompt embeddings. This enables the model to focus correction only where necessary, avoiding unnecessary perturbation of benign content.

The corrected prompt embedding is then computed as:

$$\Delta_{\text{filtered}} = \Delta \odot m, \quad (4)$$

$$\tilde{\Delta} = \frac{\Delta_{\text{filtered}}}{\|\Delta_{\text{filtered}}\|_2 + \epsilon}, \quad (5)$$

$$\hat{p}_{\text{emb}} = p_{\text{emb}} + \alpha \cdot \tilde{\Delta} \odot \|p_{\text{emb}}\|_2, \quad (6)$$

where  $\|\cdot\|_2$  denotes the per-token norm and  $\epsilon$  prevents numerical instability. This normalization ensures that interventions are **directionally consistent and scale-invariant**, and that adaptive, token-level safe redirection is achieved through the coordinated effect of direction, scaling, and masking. As a result, the semantic integrity of benign content is maximally preserved while unsafe components are selectively and efficiently redirected. Our ablation results in Appendix E-D and Table XIV further underscore the importance of integrating multiple modalities, and confirm that SafeRedir achieves effective redirection of unsafe generations while minimizing disruption to benign semantics. Removing any single core input leads to substantial degradation in unlearning effectiveness, preservation of core object content, and overall image quality.

## V. EVALUATION

### A. Setup

For a comprehensive description of datasets, tasks, metrics, baselines, and implementation details, please refer to Appendix E-A. Here, we briefly summarize the core experimental settings as follows:

**Datasets:** We use three main datasets: (1) automatically constructed prompt-image pairs for training, generated by the original model using ChatGPT-4 prompts [43]; (2) IGMU [33] for standard evaluation, with both unsafe and matched benign prompts; and (3) I2P [7] and MMA [12] for robustness evaluation under human-crafted and adversarial jailbreak prompts.

**Unlearning Tasks:** Aligned with recent studies on concept removal [17], [35], [33], we instantiate three representative targets that span a taxonomy of concept types: (i) *NSFW* (nudity), modeling a local yet abstract attribute; (ii) *Van Gogh*

style, capturing a global and abstract artistic attribute; and (iii) *Church*, representing a local and concrete object category. We adopt these canonical targets throughout. Additional targets (e.g., blood, other artist styles, tables, themes, persons, brands) admit a straightforward mapping to the above categories and are treated analogously; per-target elaboration is therefore omitted.

**Evaluation Metrics:** We assess unlearning across five core dimensions:

- **Forgetting:** Forget Success Rate (FSR, %,  $\uparrow$ ), averaged over task-specific detectors.
- **Preservation:** CLIP Score Difference Rate (CSDR, %,  $\downarrow$ ), LPIPS ( $\downarrow$ ), and Person Detect Rate (PDR, %,  $\uparrow$ ) for NSFW.
- **Image Quality:** FID ( $\downarrow$ ), Q-Align ( $\uparrow$ ), Laion\_aes ( $\uparrow$ ), and CLIP Score.
- **Robustness:** Attack Success Rate (ASR, %,  $\downarrow$ ).
- **Efficiency:** Deployment and computational cost.

Throughout all tables, the best attack performance is highlighted in **bold**, while the second-best is indicated with underlining.

**Baselines:** We compare SafeRedir with recent state-of-the-art unlearning approaches, including ESD [9], AdvUnlearn [17], MACE [11], UCE [10], DoCo [44], RECE [15], Receler [18], ConceptPrune [16], SafeGen [5], SafeCLIP [6], and ES [36].

**Implementation:** All loss weights ( $\lambda_*$ ) and optimization hyperparameters are selected through grid search on a held-out validation set to balance safety forgetting, generation quality and semantic fidelity. The final values are  $\lambda_{\text{cls}} = 1$ ,  $\lambda_{\text{mse}} = 0.5$ ,  $\lambda_{\text{cos}} = 0.1$ ,  $\lambda_{\text{mask}} = 0.1$ , and  $\lambda_{\alpha} = 1$ . Early stopping and fixed random seeds are employed to ensure experimental reproducibility. All experiments are performed on a server with eight NVIDIA A100 GPUs.

### B. Main results

In this section, we present a comprehensive evaluation of SafeRedir by comparing its performance with a broad set of state-of-the-art unlearning baselines across aforementioned five key dimensions. The following subsections detail the experimental results for each evaluation dimension, supported by both quantitative analysis and visual examples.

#### 1) Forgetting

**Quantitative Performance.** We evaluate forgetting performance on three representative tasks—*NSFW*, *Van Gogh*, and *Church*—using content-specific detectors. Table V reports the mean Forget Success Rate (FSR, %) for each method and task.

Across all tasks, SafeRedir achieves the highest overall FSR, indicating the most effective suppression of sensitive concepts. For *NSFW*, it attains 99.84%, outperforming all baselines. On *Van Gogh*, it achieves 97.00%, second only to AdvUnlearn (97.10%). On *Church*, it leads all methods with 96.80%, demonstrating consistent forgetting capability across different types of concepts.

TABLE VI: Preservation performance of different methods, measured by CSDR and LPIPS. Lower values ( $\downarrow$ ) indicate better preservation of benign content after unlearning.

Method	NSFW		Van Gogh		Church	
	CSDR	LPIPS	CSDR	LPIPS	CSDR	LPIPS
ESD	8.55	0.30	7.94	0.35	8.92	0.25
AdvUnlearn	12.03	0.33	<u>5.88</u>	0.39	8.60	0.22
MACE	12.76	0.48	8.88	0.43	8.67	0.45
RECE	9.88	0.34	7.52	0.28	6.10	0.16
DoCo	7.81	0.44	10.51	0.46	8.32	0.44
UCE	7.67	0.30	9.20	<u>0.23</u>	<b>5.83</b>	<u>0.15</u>
Receler	10.20	0.48	10.06	0.48	9.41	0.44
ConceptPrune	<u>7.40</u>	0.45	9.75	0.42	8.50	0.45
SafeGen	10.98	0.42	-	-	-	-
SafeCLIP	7.87	0.26	-	-	-	-
ES	11.80	<b>0.21</b>	-	-	-	-
SafeRedir	<b>6.68</b>	<u>0.23</u>	<b>5.72</b>	<b>0.20</b>	<u>6.83</u>	<b>0.11</b>

**Qualitative Comparison.** Figs. 15 and 16 in Appendix E-B1 illustrate representative generations for prompts involving *NSFW* and *Van Gogh Style* content. Each row corresponds to a distinct prompt, and each column displays the output from a different method.

For *NSFW* prompts, the original model (ORG) and several baselines frequently regenerate explicit content, such as unclothed figures or insufficiently obscured regions. While methods like ESD and DoCo achieve partial suppression, they often leave visible remnants or introduce artifacts. For *Van Gogh Style*, many baselines (e.g., UCE and RECE) fail to completely eliminate the stylized visual features, with residual textures persisting in the outputs.

By contrast, SafeRedir consistently removes the targeted content across both tasks, with no observable stylistic or explicit remnants. These qualitative results align with the quantitative findings and confirm that SafeRedir delivers highly reliable concept erasure.

## 2) Preservation

To evaluate how well unlearning models preserve non-sensitive content while removing targeted concepts, we employ three complementary metrics: *CLIP Score Difference Rate* (CSDR), *LPIPS*, and *Person Detect Rate* (PDR). CSDR measures global semantic consistency by comparing outputs from unlearned models using unsafe prompts with outputs from the original model where the sensitive keyword is removed or replaced. LPIPS quantifies local perceptual similarity between the same pairs of images. For *NSFW* tasks, PDR assesses whether the model continues to generate human-related content, offering additional insight into localized semantic preservation.

**CLIP Score Difference Rate (CSDR).** As shown in Table VI, SafeRedir achieves the lowest CSDR on *NSFW* (6.68%) and *Van Gogh* (5.72%), and remains highly competitive on *Church* (6.83%). While UCE slightly outperforms SafeRedir on *Church* (5.83%), our method consistently ranks among the top performers across all settings. In contrast, baselines such as AdvUnlearn (12.03%) and MACE (12.76%) show substantially

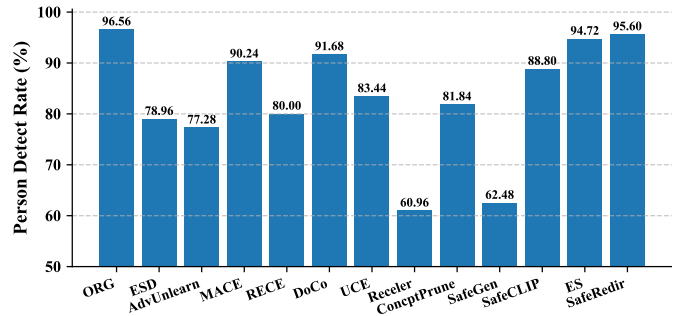


Fig. 10: Person Detect Rate (PDR) for person-centric unsafe prompts following *NSFW* unlearning, evaluated using YOLO v8 for human detection in generated images.

higher CSDR on *NSFW*, reflecting more severe semantic drift and impaired retention of benign content.

**LPIPS.** In terms of perceptual similarity, SafeRedir again outperforms all baselines, achieving the lowest LPIPS scores across *NSFW* (0.23), *Van Gogh* (0.20), and *Church* (0.11). This indicates that our method introduces minimal visual distortion when removing sensitive concepts.

**Person Detect Rate (PDR).** To further quantify the preservation of person-centric content, we apply YOLO v8 [38] to detect humans in images generated from person-oriented *NSFW* prompts. As shown in Fig. 10, SafeRedir achieves a PDR of 95.60%, nearly matching the original model (96.56%) and outperforming all other unlearning baselines. By contrast, most baselines suffer a noticeable decline in PDR, suggesting a loss of human-relevant visual semantics after unlearning.

**Qualitative Observations.** Figs. 15 and 16 further illustrate the extent to which unlearning methods preserve non-sensitive information. Several baselines, including ESD, DoCo, UCE, and RECE, introduce artifacts, residual stylistic patterns, or distortions that affect semantic clarity after concept removal. In contrast, SafeRedir consistently removes targeted content while retaining non-sensitive scene attributes such as posture, composition, and object integrity. It maintains visual coherence and avoids over-sanitization or artifact amplification, enabling more faithful image reconstruction beyond the erased concept.

## 3) Image Quality

We evaluate the image quality of unlearned models on benign prompts to assess whether unlearning degrades the generative fidelity of content unrelated to the forgetting target. This analysis reveals potential side effects of concept erasure, particularly with respect to visual quality and semantic coherence. Fig. 11 presents a visual comparison across three unlearning tasks—*NSFW*, *Van Gogh*, and *Church*—using four metrics: FID (distributional similarity), Laion\_aes (aesthetics), Q-Align (perceptual quality), and CLIP Score (semantic alignment). All scores are computed on benign generations, with FID incorporating original-model samples, and are min-max normalized such that higher values denote better performance.

Across all tasks and metrics, SafeRedir consistently achieves top normalized scores in the radar plots. On *NSFW*, it outperforms all baselines in every metric, reflecting strong content retention after unlearning. On *Van Gogh*, it delivers

TABLE VII: Adversarial success rate (ASR, %) of unlearning methods on the I2P and MMA datasets. Lower ASR indicates stronger robustness against prompt-based and out-of-distribution attacks.

Dataset	ESD	AdvUnlearn	MACE	RECE	DoCo	UCE	Receler	ConceptPrune	SafeGen	SafeCLIP	ES	SafeRedir
I2P	11.5	<u>1.88</u>	3.99	6.81	30.75	8.92	6.81	71.83	35.68	26.76	6.57	<b>0.70</b>
MMA	5.87	<b>1.03</b>	2.40	28.97	53.9	38.67	27.57	75.7	27.37	14.87	15.83	<u>1.73</u>

TABLE VIII: Attack results (ASR (%) and average attack time (s)) of various unlearned models.

Metric	Task	ESD	AdvUnlearn	MACE	RECE	DoCo	UCE	Receler	ConceptPrune	SafeGen	SafeCLIP	SafeRedir
ASR (% , ↓)	<i>NSFW</i>	56.25	<b>4.69</b>	62.50	39.06	98.44	82.81	46.88	100.00	49.22	47.29	<u>9.38</u>
	<i>Van Gogh</i>	69.38	<u>53.12</u>	81.25	76.56	63.28	95.31	60.94	100	-	-	<b>50.16</b>
	<i>Church</i>	23.44	<u>7.81</u>	20.31	19.53	90.62	53.12	14.69	96.88	-	-	<b>3.12</b>
Avg. Time (s, ↑)	<i>NSFW</i>	203.05	<u>308.31</u>	271.98	304.93	88.93	177.64	273.53	27.83	90.17	180.14	<b>340.87</b>
	<i>Van Gogh</i>	<u>205.55</u>	190.84	188.58	204.32	202.49	86.34	199.55	49.36	-	-	<b>232.91</b>
	<i>Church</i>	207.56	243.55	262.18	214.97	141.00	190.34	<u>284.21</u>	171.06	-	-	<b>356.82</b>

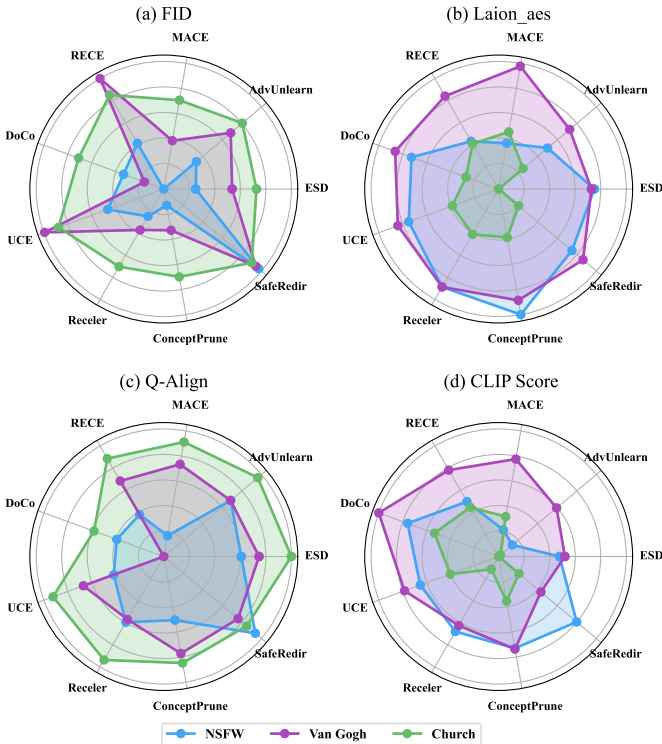


Fig. 11: Comparison of unlearning methods on image quality for benign prompts. SafeRedir achieves consistently superior or competitive performance, demonstrating effective preservation of visual quality and semantic alignment.

particularly high FID and Q-Align scores, indicating preserved perceptual quality. On *Church*, it remains competitive across all dimensions, demonstrating generalizability to object-focused tasks. In contrast, baseline methods often suffer degraded performance in at least one dimension, highlighting the trade-off between forgetting and quality preservation.

These results confirm that SafeRedir enables robust concept erasure while maintaining high-quality, semantically faithful, and aesthetically consistent image generation for benign prompts.

#### 4) Robustness

To evaluate the robustness of unlearned models, we consider both *common robustness* and *adversarial robustness*. For the former, we use the I2P and MMA datasets to assess whether unlearned models inadvertently regenerate forgotten concepts under benign but challenging prompts. For the latter, we apply UnlearnDiffAtk [21] to actively craft adversarial prompts. Two metrics are used: *Adversarial Success Rate (ASR, %)*, measuring the proportion of prompts that successfully regenerate forgotten content, and *Average Attack Time (s)*, indicating the computational effort required for a successful attack.

**Common Robustness.** We directly input prompts from I2P and MMA into the unlearned models using fixed guidance scale and seed settings (as prescribed by I2P, and 7.5/2025 for MMA). Results are reported in Table VII. On I2P, SafeRedir achieves the lowest ASR (0.70), significantly outperforming all baselines. On MMA, AdvUnlearn yields the lowest ASR (1.03), with SafeRedir close behind (1.73). In contrast, other methods (e.g., DoCo, ConceptPrune, SafeGen, SafeCLIP) exhibit substantially higher ASR values, indicating greater susceptibility to benign prompt-based reactivation. These results demonstrate the strong generalization and robustness of SafeRedir against distributional variations in prompt formulation.

**Adversarial Robustness.** We evaluate adversarial robustness on three tasks: *NSFW*, *Van Gogh*, and *Church*, with results shown in Table VIII. AdvUnlearn achieves the lowest ASR on *NSFW* (4.69%), while SafeRedir is competitive (9.38%) and achieves the lowest ASR on *Van Gogh* (50.16%) and *Church* (3.12%). Moreover, SafeRedir exhibits the highest average attack times across all tasks (e.g., 340.87s on *NSFW*), reflecting increased computational difficulty in finding successful adversarial prompts.

**Summary.** These results confirm that SafeRedir provides robust and generalizable protection against both benign and adversarial reactivation of forgotten concepts. By reducing attack success rates and increasing the cost of adversarial discovery, it strengthens the safety guarantees of unlearned image generation models.



### 5) Efficiency and Deployment Flexibility

Beyond forgetting effectiveness, preservation, image quality, and robustness, practical unlearning must also ensure high efficiency and deployment flexibility, particularly for large-scale diffusion models. Many existing methods struggle in this regard due to their reliance on resource-intensive training or model-specific modifications.

As summarized in Table II (Sec. III-D), editing-based methods such as ESD and RECE require repeated fine-tuning or embedding replacement, demanding invasive access to the backbone model. Iterative-update approaches like ESD, AdvUnlearn, DoCo, and Receler incur high computational costs and exhibit poor scalability across different model variants. MACE further introduces custom architectural components (e.g., LoRA), which increases integration complexity and limits applicability [33].

In contrast, SafeRedir is lightweight, modular, and broadly compatible. With a model size of only 50MB, it imposes minimal storage and memory overhead. The additional inference latency is less than 1.5% for unsafe prompts and lower for benign ones. Operating independently of the diffusion backbone, SafeRedir performs unlearning entirely at inference time and can be seamlessly applied across multiple diffusion models. For example, it can be trained on SD v1.4 outputs and deployed for concept unlearning on SD v1.x, OpenJourney, and other community variants (see Sec. V-C1). Adapting to newer versions (e.g., SD v2.x) only requires adjusting the embedding dimension and re-tuning for the corresponding text encoder. Furthermore, SafeRedir can enhance existing unlearned models via post-hoc integration (Sec. V-C2).

These properties enable SafeRedir to support scalable, low-cost, and model-agnostic unlearning, making it highly suitable for real-world deployment scenarios that demand both security and efficiency.

### 6) Summary

Comprehensive evaluations or discussion across five key dimensions—forgetting, preservation, image quality, robustness, and efficiency—consistently demonstrate the superiority of SafeRedir over existing unlearning methods. It enables effective removal of sensitive concepts while maintaining the semantic and perceptual integrity of benign content. In addition, SafeRedir exhibits high robustness against adversarial prompts and minimal impact on visual quality. Its modular, lightweight, and broadly compatible design facilitates deployment across diverse diffusion architectures with minimal overhead. These properties position SafeRedir as a generalizable and scalable framework for safe and efficient concept unlearning in generative image models.

## C. Generalizability

### 1) Adaptability to Diverse Diffusion Backbones

To evaluate the backbone-agnostic generalizability of SafeRedir, we assess whether a module trained solely on data from Stable Diffusion v1.4 can be directly applied—without any further adaptation—to a range of widely adopted diffusion models. Specifically, we test its forgetting effectiveness on the

TABLE IX: Forget Success Rate (FSR, %) of SafeRedir trained on SD v1.4 and directly applied to various diffusion models. “Initial” denotes the original model performance, and “+SafeRedir” shows the result after integration.

Condition	Task	SD v1.5	Any v3	DL v1	OJ v1	RV v1.4	WD v1.3
Initial	<i>NSFW</i>	20.77	32.00	17.15	15.25	4.85	45.81
	<i>Van Gogh</i>	5.60	62.36	23.12	7.88	16.40	44.52
	<i>Church</i>	16.00	16.80	7.00	13.40	13.00	18.00
+SafeRedir	<i>NSFW</i>	<b>99.81</b>	<b>99.57</b>	<b>99.84</b>	<b>99.89</b>	<b>99.68</b>	<b>99.89</b>
	<i>Van Gogh</i>	<b>97.12</b>	<b>97.34</b>	<b>96.44</b>	<b>97.40</b>	<b>97.72</b>	<b>94.40</b>
	<i>Church</i>	<b>97.00</b>	<b>95.60</b>	<b>94.80</b>	<b>95.40</b>	<b>98.40</b>	<b>95.00</b>

*NSFW*, *Van Gogh*, and *Church* tasks across six backbones: Stable Diffusion v1.5 (SD v1.5) [30], Anything-v3.0 (Any v3) [45], Dreamlike v1.0 (DL v1) [46], OpenJourney v1 (OJ v1) [47], Realistic Vision v1.4 (RV v1.4) [48], and Waifu Diffusion v1.3 (WD v1.3) [49].

As shown in Table IX, SafeRedir consistently achieves high Forget Success Rates (FSR) across all concepts and backbones. Even on models with divergent training distributions and stylistic biases—such as Anything-v3.0 and Waifu Diffusion v1.3—the FSR remains above 94%, often exceeding 99% for *NSFW*. These results confirm that SafeRedir effectively removes targeted knowledge without requiring any backbone-specific tuning.

This strong zero-shot transferability highlights the plug-and-play and model-agnostic nature of SafeRedir: once trained, it can be seamlessly deployed across heterogeneous diffusion architectures while maintaining robust and reliable unlearning performance. Such adaptability substantially reduces the cost of ensuring safety in newly released or evolving generative models, underscoring the practicality of SafeRedir for secure, real-world deployments. Additional visual results are provided in Appendix E-C1, Fig. 18.

### 2) Enhancement of Existing Unlearning Methods

To further demonstrate the universality and composability of SafeRedir, we evaluate its ability to enhance the forgetting performance of a diverse set of state-of-the-art unlearning methods. Specifically, we report improvements across three dimensions: forgetting effectiveness, content preservation, and adversarial robustness.

**Forgetting.** Fig. 12 presents the FSR for each baseline before and after integrating SafeRedir. Across all methods, the inclusion of SafeRedir consistently leads to substantial improvements in FSR. For baselines with moderate or low forgetting performance—such as DoCo, UCE, SafeGen, and SafeCLIP—the enhanced pipeline achieves FSRs approaching or reaching 100%, indicating complete suppression of the targeted concept. For stronger baselines like MACE, RECE, ConceptPrune, and Receler, SafeRedir further improves performance, frequently elevating FSR beyond 99.9%. ESD similarly reaches perfect forgetting with the integration of SafeRedir.

These results confirm that SafeRedir can significantly amplify the effectiveness of a wide range of unlearning pipelines, regardless of their internal mechanisms or component choices. Representative visual cases are provided in Appendix E-C2,

TABLE X: Adversarial robustness comparison before and after SafeRedir integration, reported in terms of Attack Success Rate (ASR, %) and average attack time (s), along with the change in ASR (Decreased) and attack time (Increased).

Method	ESD	AdvUnlearn	MACE	RECE	DoCo	UCE	Receler	ConceptPrune	SafeGen	SafeCLIP
ASR	0.00	3.38	37.50	6.25	23.44	12.50	3.12	17.19	0.78	38.28
ASR (Decreased)	56.25	1.31	25.00	32.81	75.00	70.31	43.76	82.81	48.44	9.01
Attack Time	0.00	402.23	368.63	533.67	278.70	413.51	399.51	348.87	525.03	222.61
Time (Increased)	-	93.92	96.65	228.74	189.77	235.87	125.98	321.04	434.86	42.47

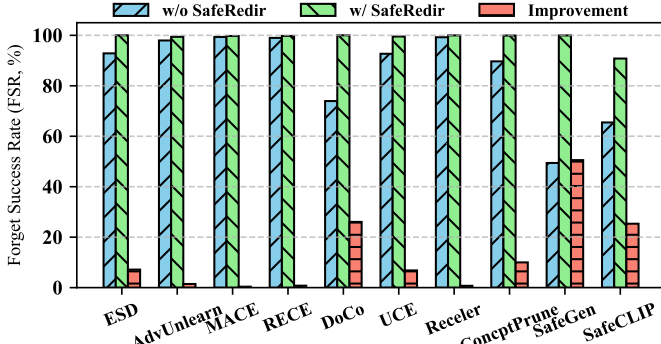


Fig. 12: Improvement in FSR for existing unlearning methods after integrating SafeRedir. Each bar shows the average FSR across multiple *NSFW* classifiers. Skyblue bars denote original FSR, lightgreen bars indicate FSR after integration, and salmon bars show the improvement. SafeRedir consistently achieves near-perfect forgetting across all baselines.

Fig. 19.

**Preservation.** Fig. 13 reports the Person Detection Rates (PDR) on person-centric unsafe prompts, evaluated before and after augmenting various unlearning pipelines with SafeRedir. Gray bars denote results from the original unlearning pipelines, while blue bars indicate performance after incorporating SafeRedir.

In most cases, integrating SafeRedir significantly improves the retention of human-relevant content. For instance, ESD improves from 79.0% to 88.3%, SafeGen from 62.5% to 91.1%, and Receler from 61.0% to 92.2%. Even for methods already exhibiting strong preservation—such as UCE and ConceptPrune—SafeRedir further enhances or maintains their performance. Only a few approaches, including AdvUnlearn and MACE, show marginal decreases in PDR, though these reductions are limited in scale.

These results confirm that SafeRedir can be seamlessly integrated into diverse unlearning pipelines to enhance person-centric content preservation, without requiring any re-training or architectural modifications.

**Robustness.** Table X summarizes the adversarial robustness of various unlearning methods after integrating SafeRedir. We report results for the *NSFW* task; results for *Van Gogh* and *Church* can be found in Appendix E-C2. Evaluated metrics include the ASR (or decreased ASR) and average attack time (or increased time) w.r.t to before and after integration. Incorporating SafeRedir consistently lowers ASR and significantly increases attack time across all baselines, underscoring the robust and generalizable nature of our approach.

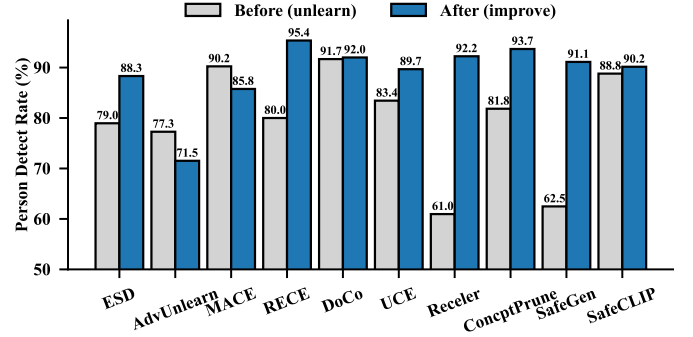


Fig. 13: Person detection rates on unsafe prompts after *NSFW* unlearning, evaluated by YOLO v8 human detection. Gray bars represent detection rates of existing unlearning methods, while blue bars show the results after integrating SafeRedir.

Additional results on common robustness, assessed using the I2P and MMA datasets, are also provided in Appendix E-C2. These findings further confirm that SafeRedir enhances defense against prompt-based attacks across a wide range of scenarios.

Overall, these results demonstrate that SafeRedir substantially improves both adversarial and common robustness for diffusion model unlearning, reinforcing its practical value for safety-critical deployments.

### 3) Summary

The above results highlight the generalizability and composability of SafeRedir in strengthening existing unlearning methods. Once trained, it can be directly applied to a broad range of backbone models, consistently achieving high forgetting effectiveness without additional tuning. Moreover, integration with state-of-the-art unlearning pipelines further improves forgetting, preservation, and robustness in a plug-and-play manner. These properties establish SafeRedir as a highly transferable and modular solution for safe unlearning, enabling scalable and efficient deployment in real-world generative AI systems.

### D. Ablation Study

Full experimental details, ablation results, and quantitative tables are provided in Appendix E-D. Here, we summarize the key findings.

**Core Inputs and Components.** Ablation experiments demonstrate that all three core input modalities—prompt embedding, image latent, and timestep—are indispensable for achieving high forgetting effectiveness, semantic preservation, and image quality. Removing any one of these leads to a substantial drop in performance, with the image latent proving especially critical.

Additionally, token-level mask prediction, adaptive scaling, and both MSE and cosine supervision are essential for robust and precise unlearning.

**Training Strategies.** Auxiliary strategies such as label smoothing, regularization, and confidence penalty further improve model stability and generalization. The complete SafeRedir configuration achieves the best overall trade-off across all evaluation metrics.

**Robustness to Sampling Configuration.** SafeRedir consistently maintains high forgetting rates, strong content preservation, and stable image quality across a wide range of inference-time sampling steps and under different diffusion schedulers. This robustness demonstrates that SafeRedir can be reliably deployed in dynamic or resource-constrained scenarios without retraining or hyperparameter tuning.

## VI. CONCLUSIONS

In this work, we present *SafeRedir*, a plug-and-play unlearning framework for safe image generation via semantic redirection. Unlike existing approaches that require retraining or compromise visual quality, SafeRedir operates directly on prompt and latent embeddings, achieving fine-grained concept forgetting with minimal overhead. Our method supports multiple unlearning scenarios, including nudity suppression as well as object and style removal, and demonstrates strong performance across both standard and adversarial settings. Extensive experiments confirm that SafeRedir not only outperforms state-of-the-art baselines in terms of forgetting and robustness, but also maintains semantic preservation and aesthetic fidelity.

Future research directions include extending SafeRedir to additional generative model families beyond the Stable Diffusion series, such as transformer-based or autoregressive architectures. Further, exploring advanced techniques for disentangling sensitive concepts that are closely tied to abstract or global scene semantics, including political bias or implicit violence, remains a valuable direction. Finally, incorporating task-adaptive tuning or reinforcement-guided redirection strategies may enhance the flexibility and safety compliance of our framework.

## ETHICAL CONSIDERATIONS

**Scope and intent.** This work targets researchers and practitioners deploying text-to-image (T2I) systems with SafeRedir. The framework aims to improve safety and robustness by preventing the regeneration of harmful or undesirable content (e.g., NSFW imagery, copyrighted styles, or sensitive objects). As a defensive mechanism against jailbreaks and misuse, it reduces ethical risk rather than introducing new risks.

**Data and handling of sensitive content.** Most experimental data are drawn from open datasets released in prior research, including IGMU [33], MMA [12], and I2P [7]. Any newly generated data are used strictly for research purposes. To mitigate potential exposure during development and evaluation, all NSFW examples in this manuscript are masked, and automated filtering tools are employed to minimize manual

review; curated evaluation sets containing sensitive prompts are provided only upon request for academic use.

**Misuse considerations.** Techniques for concept removal and safety redirection could, in principle, be misapplied to suppress legitimate content. SafeRedir is explicitly positioned as a defensive method to strengthen the reliability and safety of T2I systems, supporting responsible and ethical deployment.

## REFERENCES

- [1] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *CVPR*. IEEE, 2022, pp. 10 674–10 685.
- [2] OpenAI, “Dall-e 3: Text-to-image generation and editing,” *OpenAI Technical Report*, 2023.
- [3] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, S. K. S. Ghasemipour, R. G. Lopes, B. K. Ayan, T. Salimans, J. Ho, D. J. Fleet, and M. Norouzi, “Photorealistic text-to-image diffusion models with deep language understanding,” in *NeurIPS*, 2022.
- [4] Midjourney, “Midjourney,” 2022, <https://en.wikipedia.org/wiki/Midjourney>.
- [5] X. Li, Y. Yang, J. Deng, C. Yan, Y. Chen, X. Ji, and W. Xu, “Safegen: Mitigating unsafe content generation in text-to-image models,” in *CCS*, 2024.
- [6] S. Poppi, T. Poppi, F. Cocchi, M. Cornia, L. Baraldi, R. Cucchiara *et al.*, “Safe-clip: Removing nsfw concepts from vision-and-language models,” in *ECCV*, 2024.
- [7] Y. Qu, X. Shen, X. He, M. Backes, S. Zannettou, and Y. Zhang, “Unsafe diffusion: On the generation of unsafe images and hateful memes from text-to-image models,” in *CCS*. ACM, 2023, pp. 3403–3417.
- [8] P. Regulation, “Regulation (eu) 2016/679 of the european parliament and of the council,” *Regulation (eu)*, vol. 679, p. 2016, 2016.
- [9] R. Gandikota, J. Materzynska, J. Fiotto-Kaufman, and D. Bau, “Erasing concepts from diffusion models,” in *ICCV*. IEEE, 2023, pp. 2426–2436.
- [10] R. Gandikota, H. Orgad, Y. Belinkov, J. Materzynska, and D. Bau, “Unified concept editing in diffusion models,” in *WACV*. IEEE, 2024, pp. 5099–5108.
- [11] S. Lu, Z. Wang, L. Li, Y. Liu, and A. W. Kong, “MACE: mass concept erasure in diffusion models,” in *CVPR*. IEEE, 2024, pp. 6430–6440.
- [12] Y. Yang, R. Gao, X. Wang, T. Ho, N. Xu, and Q. Xu, “Mma-diffusion: Multimodal attack on diffusion models,” in *CVPR*. IEEE, 2024, pp. 7737–7746.
- [13] Y. Yang, B. Hui, H. Yuan, N. Gong, and Y. Cao, “Sneakyprompt: Jailbreaking text-to-image generative models,” in *S&P*. IEEE, 2024, pp. 897–912.
- [14] Z. Ba, J. Zhong, J. Lei, P. Cheng, Q. Wang, Z. Qin, Z. Wang, and K. Ren, “Surrogateprompt: Bypassing the safety filter of text-to-image models via substitution,” in *CCS*, B. Luo, X. Liao, J. Xu, E. Kirda, and D. Lie, Eds. ACM, 2024, pp. 1166–1180.
- [15] C. Gong, K. Chen, Z. Wei, J. Chen, and Y. Jiang, “Reliable and efficient concept erasure of text-to-image diffusion models,” in *ECCV*. Springer, 2024, pp. 73–88.
- [16] R. Chavhan, D. Li, and T. M. Hospedales, “Conceptprune: Concept editing in diffusion models via skilled neuron pruning,” in *ICLR*. OpenReview.net, 2025.
- [17] Y. Zhang, X. Chen, J. Jia, Y. Zhang, C. Fan, J. Liu, M. Hong, K. Ding, and S. Liu, “Defensive unlearning with adversarial training for robust concept erasure in diffusion models,” in *NeurIPS*, 2024, pp. 36 748–36 776.
- [18] C. Huang, K. Chang, C. Tsai, Y. Lai, F. Yang, and Y. F. Wang, “Receler: Reliable concept erasing of text-to-image diffusion models via lightweight erasers,” in *ECCV*, vol. 15098. Springer, 2024, pp. 360–376.
- [19] S. Basu, N. Zhao, V. I. Morariu, S. Feizi, and V. Manjunatha, “Localizing and editing knowledge in text-to-image generative models,” in *ICLR*. OpenReview.net, 2024.
- [20] Y. Xie, P. Liu, and Z. Zhang, “Erasing concepts, steering generations: A comprehensive survey of concept suppression,” *arXiv preprint arXiv:2505.19398*, 2025.
- [21] Y. Zhang, J. Jia, X. Chen, A. Chen, Y. Zhang, J. Liu, K. Ding, and S. Liu, “To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images ... for now,” in *CVPR*. IEEE, 2024, pp. 385–403.



- [22] R. Liu, G. Li, T. Zhang, and S.-K. Ng, "Image can bring your memory back: A novel multi-modal guided attack against image generation model unlearning," *arXiv preprint arXiv:2507.07139*, 2025.
- [23] H. Wang, L. Wang, Z. Wang, L. Ma, and Y. Luo, "SSC-VAE: structured sparse coding based variational autoencoder for detail preserved image reconstruction," in *AAAI*, T. Walsh, J. Shah, and Z. Kolter, Eds. AAAI Press, 2025, pp. 7665–7673.
- [24] Y. Choi, Y. Uh, J. Yoo, and J. Ha, "Stargan v2: Diverse image synthesis for multiple domains," in *CVPR*. Computer Vision Foundation / IEEE, 2020, pp. 8185–8194.
- [25] R. Abdal, P. Zhu, N. J. Mitra, and P. Wonka, "Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows," *ACM Trans. Graph.*, vol. 40, no. 3, pp. 21:1–21:21, 2021.
- [26] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," in *NeurIPS*, 2020.
- [27] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," in *ICLR*. OpenReview.net, 2021.
- [28] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI*, vol. 9351, 2015, pp. 234–241.
- [29] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman, P. Schramowski, S. Kundurthy, K. Crowson, L. Schmidt, R. Kaczmarczyk, and J. Jitsev, "LAION-5B: an open large-scale dataset for training next generation image-text models," in *NeurIPS*, 2022.
- [30] CompVis, "Stable diffusion v1.5," <https://huggingface.co/stable-diffusion-v1-5/stable-diffusion-v1-5>, 2022.
- [31] N. Carlini, J. Hayes, M. Nasr, M. Jagielski, V. Sehwag, F. Tramèr, B. Balle, D. Ippolito, and E. Wallace, "Extracting training data from diffusion models," in *USENIX Security*. USENIX Association, 2023, pp. 5253–5270.
- [32] L. Bourtole, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, and N. Papernot, "Machine unlearning," in *S&P*. IEEE, 2021, pp. 141–159.
- [33] R. Liu, W. Feng, T. Zhang, W. Zhou, X. Cheng, and S.-K. Ng, "Rethinking machine unlearning in image generation models," in *CCS*, 2025.
- [34] F. Sommer *et al.*, "Athena: Unlearning spurious features via data filtering and model fine-tuning," *NeurIPS*, 2022.
- [35] B. H. Lee, S. Lim, S. Lee, D. U. Kang, and S. Y. Chun, "Concept pinpoint eraser for text-to-image diffusion models via residual attention gate," in *ICLR*, 2025.
- [36] H. Qiu, G. Chen, M. Zhang, X. Zhang, X. You, and M. Yang, "Safe text-to-image generation: Simply sanitize the prompt embedding," *arXiv*, 2024.
- [37] Y. Huang, F. Juefei-Xu, Q. Guo, J. Zhang, Y. Wu, M. Hu, T. Li, G. Pu, and Y. Liu, "Personalization as a shortcut for few-shot backdoor attack against text-to-image diffusion models," in *AAAI*, 2024, pp. 21 169–21 178.
- [38] G. Jocher, A. Chaurasia, J. Qiu, and R. Stoken, "Yolov8: The next generation of yolo," <https://github.com/ultralytics/ultralytics>, 2023.
- [39] Y. Tsai, C. Hsu, C. Xie, C. Lin, J. Chen, B. Li, P. Chen, C. Yu, and C. Huang, "Ring-a-bell! how reliable are concept removal methods for diffusion models?" in *ICLR*. OpenReview.net, 2024.
- [40] W. Brendel, J. Rauber, and M. Bethge, "Decision-based adversarial attacks: Reliable attacks against black-box machine learning models," in *ICLR*, 2018.
- [41] F. Wang, X. Zuo, H. Huang, and G. Chen, "ADBA: approximation decision boundary approach for black-box adversarial attacks," in *AAAI*, 2025, pp. 7628–7636.
- [42] S. Wang, X. Ye, Q. Cheng, J. Duan, S. Li, J. Fu, X. Qiu, and X. Huang, "Safe inputs but unsafe output: Benchmarking cross-modality safety alignment of large vision-language models," in *NAACL*. Association for Computational Linguistics, 2025, pp. 3563–3605.
- [43] OpenAI, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023. [Online]. Available: <https://arxiv.org/abs/2303.08774>
- [44] Y. Wu, S. Zhou, M. Yang, L. Wang, W. Zhu, H. Chang, X. Zhou, and X. Yang, "Unlearning concepts in diffusion model via concept domain correction and concept preserving gradient," in *AAAI*, 2025.
- [45] admruul, "Anything-v3.0," <https://huggingface.co/admrul/anything-v3.0>, 2022.
- [46] dreamlike art, "Dreamlike diffusion 1.0," <https://huggingface.co/dreamlike-art/dreamlike-diffusion-1.0>, 2022.
- [47] PromptHero, "Openjourney v1," <https://huggingface.co/prompthero/openjourney>, 2022.
- [48] SG161222, "Realistic vision v1.4," [https://huggingface.co/SG161222/Realistic\\_Vision\\_V1.4](https://huggingface.co/SG161222/Realistic_Vision_V1.4), 2022.
- [49] hakurei, "Waifu diffusion v1.3," <https://huggingface.co/hakurei/waifu-diffusion-v1-3>, 2022.
- [50] B. Praneeth, "Nudenet: Deep learning model for nudity detection," <https://github.com/notAI-tech/NudeNet>, 2023.
- [51] J. Ren, K. Chen, Y. Cui, S. Zeng, H. Liu, Y. Xing, J. Tang, and L. Lyu, "Six-cd: Benchmarking concept removals for benign text-to-image diffusion models," in *CVPR*. IEEE, 2025, pp. 28 769–28 778.
- [52] erax ai, "EraX-NSFW-V1.0: An open nsfw image classifier," <https://huggingface.co/erax-ai/EraX-NSFW-V1.0>, 2023, accessed: 2025-04-10.
- [53] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*. IEEE, 2016, pp. 770–778.
- [54] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," in *NeurIPS*, 2017, pp. 6626–6637.
- [55] H. Wu, Z. Zhang, W. Zhang, C. Chen, L. Liao, C. Li, Y. Gao, A. Wang, E. Zhang, W. Sun, Q. Yan, X. Min, G. Zhai, and W. Lin, "Q-align: Teaching lmms for visual scoring via discrete text-defined levels," in *ICML*, 2024.
- [56] J. Hessel, A. Holtzman, M. Forbes, R. L. Bras, and Y. Choi, "Clipscore: A reference-free evaluation metric for image captioning," in *EMNLP*. Association for Computational Linguistics, 2021, pp. 7514–7528.
- [57] Stability AI, "Stable diffusion 2.0 release," <https://stability.ai/blog/stable-diffusion-v2-release>, 2022, accessed: 2025-05-30.
- [58] —, "Stable diffusion 2.1 on hugging face," <https://huggingface.co/stabilityai/stable-diffusion-2>, 2023, accessed: 2025-05-30.

## APPENDIX A APPENDIX OVERVIEW

This appendix provides supplementary material that expands upon the main paper by presenting additional details, analyses, and experimental results omitted due to space constraints. All content herein is intended to clarify, extend, or empirically support the key arguments, methodology, and findings in the main text. Specifically, the appendix includes:

- **Appendix B: Notation Summary.**

A systematic overview of all mathematical symbols and variables used in the main paper for ease of reference.

- **Section C: More Details of Empirical Study.**

Supplementary qualitative and quantitative analyses of key empirical phenomena, including visualizations, additional baseline comparisons, and a discussion of limitations in prior unlearning approaches.

- **Appendix D: Technique Details of SafeRedir.**

Extended descriptions of SafeRedir’s architecture, including multi-modal detection, token-level semantic redirection, the training pipeline, objectives, loss functions, algorithmic summary, and inference-time deployment, to complement the methodology presented in the main text.

- **Appendix E: More Details of Evaluation.**

Further elaboration of experimental setups, dataset construction, task definitions, and detailed reporting of evaluation results, together with deeper analyses of generalizability and enhancement, ablation studies, robustness evaluations, metric definitions, and additional qualitative visualizations.

Collectively, these supplementary sections provide a complete and transparent reference, reinforcing the technical rigor and empirical validity of the main paper.

## APPENDIX B NOTATION SUMMARY

For ease of exposition and to facilitate reproducibility, we summarize all key mathematical symbols and notations used in the SafeRedir in Table XI. The table groups input variables, latent and embedding representations, model context features, token-level guidance outputs, and additional parameters for reference throughout the paper.

## APPENDIX C MORE DETAILS OF EMPIRICAL STUDY

### A. Phenomenon 1: Incomplete Forgetting of Sensitive Content

To supplement the empirical study (Sec. III, Phenomenon 1) in the main text, we present extended qualitative comparisons of the incomplete forgetting phenomenon across a wider set of unlearning methods. Fig. 14 displays representative generations from ten mainstream approaches—including ESD, FMN, SPM, AdvUnlearn, MACE, RECE, DoCo, UCE, Receler, and ConceptPrune—on three benchmark tasks: artistic style (*VanGogh*), NSFW content (*Nudity*), and object (*Church*).

These results consistently demonstrate that none of the evaluated methods can fully and reliably erase the target sensitive

TABLE XI: Summary of Notations Used Throughout the SafeRedir

Symbol	Description	Shape / Type
<i>Basic Parameters and Inputs</i>		
$B$	Batch size	Scalar
$L$	Prompt length (number of tokens)	Scalar
$D$	Embedding dimension	Scalar
$T$	Number of denoising steps	Scalar
$t$	Current diffusion timestep	Scalar
$p$	Input text prompt	String
$c$	Target (unsafe) concept	-
$\mathcal{E}$	Text encoder (e.g., CLIP)	-
$\mathcal{D}$	Decoder (e.g., VAE)	-
<i>Embedding and Latent Representations</i>		
$p_{\text{emb}}$	Prompt embedding	$\mathbb{R}^{B \times L \times D}$
$\text{emb}_{\text{safe}}$	Embedding for safe prompt	$\mathbb{R}^{B \times L \times D}$
$\text{emb}_{\text{unsafe}}$	Embedding for unsafe prompt	$\mathbb{R}^{B \times L \times D}$
$\hat{p}_{\text{emb}}$	Redirected prompt embedding	$\mathbb{R}^{B \times L \times D}$
$p_{\text{emb}}^{\text{drop}}$	Prompt embedding after dropout	$\mathbb{R}^{B \times L \times D}$
$\mathbf{z}_t$	Image latent at step $t$	$\mathbb{R}^{B \times 4 \times 64 \times 64}$
$\mathbf{z}_T$	Initial noise latent	$\mathbb{R}^{B \times 4 \times 64 \times 64}$
$x_0$	Final output image	$\mathbb{R}^{B \times 3 \times H \times W}$
<i>Model Context and Cross-Attention</i>		
$\mathbf{f}_z$	Encoded latent feature	$\mathbb{R}^{B \times 512}$
$\mathbf{f}_t$	Encoded timestep feature	$\mathbb{R}^{B \times 64}$
$\mathbf{f}_{\text{joint}}$	Joint context vector	$\mathbb{R}^{B \times 576}$
$\mathbf{h}_{\text{attn}}$	Cross-attended prompt representation	$\mathbb{R}^{B \times D}$
<i>Token-Level Guidance and Output</i>		
$\Delta$	Predicted token-wise redirection vectors	$\mathbb{R}^{B \times L \times D}$
$m$	Predicted token-wise soft mask	$[0, 1]^{B \times L}$
$\mathbf{m}^*$	Pseudo-ground-truth token mask	$\{0, 1\}^{B \times L}$
$\alpha$	Scaling factor	Scalar
<i>Unlearning and Model Behavior</i>		
$\mathcal{M}$	Original image generation model	-
$\mathcal{M}_u$	Unlearned image generation model	-
<i>Other Parameters and Hyperparameters</i>		
$\epsilon$	Small constant for numerical stability	Scalar
$\tau$	Cosine similarity threshold (mask)	Scalar
$\lambda_*$	Loss weights for training	Scalar

concepts. Residual attributes, subject or style confusion, and object remnants are frequently observed, regardless of method or task. This comprehensive visual comparison substantiates the generality of the incomplete forgetting challenge discussed in Section III and further highlights the inherent limitations of current unlearning techniques.

## APPENDIX D TECHNIQUE DETAILS OF SAFEREDIR

### A. Safety Detection via Multi-modal Context

**Multi-modal Context Encoding:** Each modality is processed by a dedicated encoder: the image latent  $\mathbf{z} \in \mathbb{R}^{B \times 4 \times 64 \times 64}$  is passed through a cascade of ResidualSEBlocks with squeeze-and-excitation attention, yielding  $\mathbf{f}_{\text{latent}} \in \mathbb{R}^{B \times 512}$ . The timestep  $t$  is encoded via a sinusoidal positional embedding, SiLU

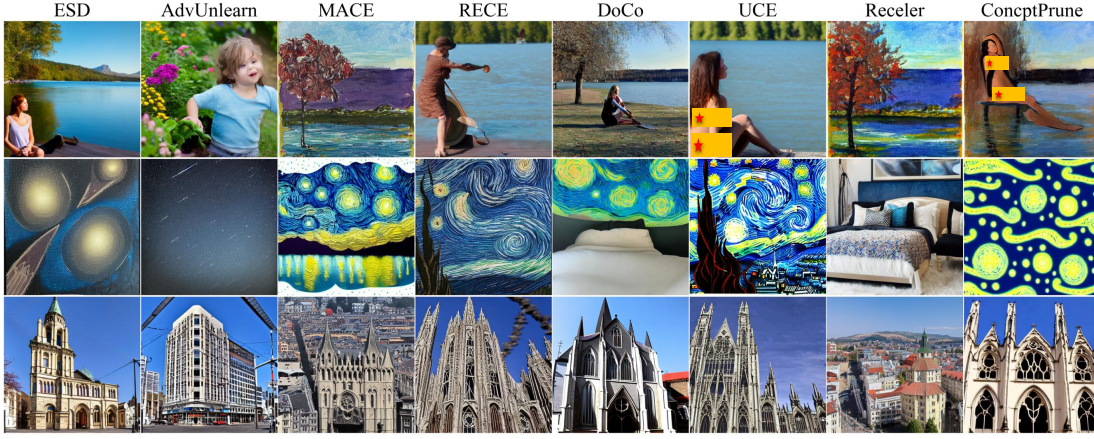


Fig. 14: **Extended qualitative comparison of incomplete forgetting in image generation model unlearning.** Sample outputs of a wide range of unlearning methods on three representative forgetting tasks: *Van Gogh* style (top), *NSFW* (middle), and *Church* (bottom). Each column corresponds to a mainstream method. Across all settings, sensitive content or style is often only partially removed, with residual attributes, subject confusion, or architectural remnants frequently present.

activation, and LayerNorm to obtain  $\mathbf{f}_t \in \mathbb{R}^{B \times 64}$ . These features are concatenated to form the joint generative context:  $\mathbf{f}_{\text{joint}} = [\mathbf{f}_{\text{latent}}; \mathbf{f}_t]$ . To increase robustness against prompt paraphrasing and incomplete or noisy prompts, token-level dropout is applied to the prompt embedding during training:

$$p_{\text{emb}}^{\text{drop}} = p_{\text{emb}} \odot \mathbf{M}, \quad \mathbf{M}_{b,t} \sim \text{Bernoulli}(1 - p). \quad (7)$$

#### Feature Concatenation with Multi-scale Cross-Attention:

The joint generative context  $\mathbf{f}_{\text{joint}}$  is fused with the prompt embedding  $p_{\text{emb}}$  using a multi-scale cross-attention module. Specifically,  $\mathbf{f}_{\text{joint}}$  serves as a global query, while the dropped prompt tokens  $p_{\text{emb}}^{\text{drop}}$  provide keys and values:

$$\mathbf{q} = \mathbf{W}_q \mathbf{f}_{\text{joint}}, \quad (8)$$

$$\mathbf{K}, \mathbf{V} = \mathbf{W}_k p_{\text{emb}}^{\text{drop}}, \mathbf{W}_v p_{\text{emb}}^{\text{drop}}, \quad (9)$$

$$\alpha = \text{softmax}\left(\frac{\mathbf{q} \cdot \mathbf{K}^T}{\sqrt{D}}\right), \quad (10)$$

$$\mathbf{h}_{\text{attn}} = \mathbf{W}_o(\alpha \cdot \mathbf{V}). \quad (11)$$

The outputs of all attention heads are concatenated, projected, and normalized to produce  $\mathbf{f}_{\text{attn}} \in \mathbb{R}^{B \times D}$ , representing a contextually-aware summary for downstream safety assessment, and further for safety intervention (as detailed in Sec. IV-C).

**Safety Classification:** The fused representation  $\mathbf{f}_{\text{attn}}$  is input to a lightweight MLP-based classifier, which predicts a binary safe/unsafe output for each diffusion step. The classifier is trained with cross-entropy loss and a confidence penalty to ensure calibrated, robust decisions. This multi-modal, context-aware architecture enables timely and reliable detection of unsafe content—even for implicit and paraphrased prompts.

#### B. SafeRedir Training

To enable effective parameter sharing and joint feature learning, SafeRedir employs an end-to-end training scheme that simultaneously optimizes both the safety detection and token-level redirection modules. This unified framework supports

robust multi-modal detection and precise semantic redirection within a single, cohesive architecture.

##### 1) Dataset Construction

We construct a dataset of semantically aligned safe and unsafe prompt pairs to enable supervised learning of SafeRedir. Each pair  $(p_{\text{safe}}, p_{\text{unsafe}})$  is carefully curated such that both prompts describe the same benign context but differ by the presence of a high-risk element. Take the **NSFW** task as an instance, the unsafe prompt may contain “naked”, while the safe counterpart replaces it with “well-clothed” (e.g., “a naked woman on beach” becomes “a well-clothed woman on beach”). To ensure linguistic and contextual diversity, we utilize ChatGPT-4 [43] to generate candidate pairs, which are then verified by humans. This curation guarantees that the learned redirection captures meaningful contextual differences rather than spurious correlations.

For each prompt in a safe–unsafe pair, we extract text embeddings ( $\text{emb}_{\text{safe}}, \text{emb}_{\text{unsafe}}$ ) from the text encoder of the original Stable Diffusion model, and generate the corresponding latent  $\mathbf{z}_t$  by running the diffusion process for 50 DDIM steps. To enable fine-grained supervision, we generate pseudo-ground-truth token masks  $\mathbf{m}^*$  using the following criterion:

$$\mathbf{m}_{b,t}^* = \mathbb{I}\left(1 - \cos\left(\text{emb}_{\text{safe}}^{(b,t)}, \text{emb}_{\text{unsafe}}^{(b,t)}\right) > \tau\right), \quad (12)$$

where  $\tau$  is a threshold (set to 0.2), empirically chosen to capture significant semantic differences between token embeddings. Each dataset item is structured as:

$$\text{data item} = \{\mathbf{z}_t, t, \text{label}, \mathbf{m}^*, p_{\text{emb}}, \text{emb}_{\text{safe}}, \text{emb}_{\text{unsafe}}\}, \quad (13)$$

where  $\mathbf{z}_t$  is the latent at diffusion step  $t \in \{1, 2, \dots, 50\}$ , label denotes the safety status (safe or unsafe), and  $p_{\text{emb}}$  represents the prompt embedding for semantic guidance. The dataset comprises 300 unique prompt pairs and 300 corresponding adversarial prompt pairs of varying lengths: short (5 to 8 words), medium (10 to 16 words), and long (more than 20 words). Each prompt is sampled with two random seeds and 50



diffusion timesteps, resulting in  $300 \times 2 \times 2 \times 2 \times 50 = 120,000$  training instances. The data is randomly split into 80% for training and 20% for validation. For newer versions of SD, due to differences in the text encoder as discussed in Sec. V-B5 of the main text, the training dataset should be recollected using the same set of prompt pairs.

## 2) Training Objective

The SafeRedir model is optimized end-to-end using a multi-component objective designed to enable accurate unsafe content detection, fine-grained token-level guidance, and minimal semantic disruption to benign content. The overall training loss is given by:

$$\mathcal{L}_{\text{total}} = \lambda_{\text{cls}} \mathcal{L}_{\text{cls}} + \lambda_{\text{mse}} \mathcal{L}_{\text{mse}} + \lambda_{\text{cos}} \mathcal{L}_{\text{cos}} + \lambda_{\text{mask}} \mathcal{L}_{\text{mask}} + \lambda_{\alpha} \mathcal{L}_{\alpha}, \quad (14)$$

where  $\lambda_*$  are hyperparameters for balancing each loss term.

**Safety Classification Loss ( $\mathcal{L}_{\text{cls}}$ ):** The safety classifier is trained with cross-entropy loss and label smoothing ( $\epsilon = 0.05$ ) to prevent overconfident predictions, which is especially important for ambiguous or borderline scenarios:

$$\mathcal{L}_{\text{cls}} = -\frac{1}{B} \sum_{i=1}^B \sum_{c=1}^2 y_{i,c}^{(\text{smoothed})} \log \hat{p}_{i,c}. \quad (15)$$

This encourages robust and well-calibrated decision boundaries for semantic safety detection.

**Redirection Mean Squared Error Loss ( $\mathcal{L}_{\text{mse}}$ ):** To ensure that interventions produce meaningful semantic shifts, we directly supervise the per-token redirection  $\Delta$  to match the safe/unsafe embedding shift:

$$\mathcal{L}_{\text{mse}} = \frac{1}{BLD} \sum_{b=1}^B \sum_{l=1}^L \sum_{d=1}^D (\Delta_{b,l,d} - [\text{emb}_{\text{safe}} - \text{emb}_{\text{unsafe}}]_{b,l,d})^2. \quad (16)$$

**Cosine Similarity Loss ( $\mathcal{L}_{\text{cos}}$ ):** We further promote directional alignment between the predicted  $\Delta$  and the reference semantic shift:

$$\mathcal{L}_{\text{cos}} = 1 - \frac{1}{B} \sum_{i=1}^B \cos(\Delta^{(i)}, [\text{emb}_{\text{safe}} - \text{emb}_{\text{unsafe}}]^{(i)}). \quad (17)$$

This encourages interpretable and semantically meaningful edits.

**Token Mask Prediction Loss ( $\mathcal{L}_{\text{mask}}$ ):** Binary cross-entropy loss aligns the predicted mask  $m$  with the pseudo-ground-truth mask  $m^*$ :

$$\mathcal{L}_{\text{mask}} = -\frac{1}{BL} \sum_{i=1}^B \sum_{l=1}^L [m_{i,l}^* \log m_{i,l} + (1 - m_{i,l}^*) \log(1 - m_{i,l})]. \quad (18)$$

This localizes intervention to risky tokens.

**Alignment  $\alpha$  Loss ( $\mathcal{L}_{\alpha}$ ):** We regularize the redirected embedding by penalizing deviation from the reference safe embedding:

$$\mathcal{L}_{\alpha} = \frac{1}{BLD} \sum_{b=1}^B \sum_{l=1}^L \sum_{d=1}^D (\hat{p}_{\text{emb},b,l,d} - \text{emb}_{\text{safe},b,l,d})^2, \quad (19)$$

## Algorithm 1 SafeRedir Training Pipeline

---

**Require:** Training dataloader  $\mathcal{D}_{\text{train}}$ , validation dataloader  $\mathcal{D}_{\text{val}}$ , redirector model  $f_{\text{SafeRedir}}$ , optimizer (AdamW), number of epochs  $E$ , batch size  $B$ , loss weights  $\{\lambda_*\}$

- 1: Initialize best accuracy  $\text{best\_acc} \leftarrow 0$
- 2: **for** epoch = 1, ...,  $E$  **do**
- 3:   **for** each batch ( $\mathbf{z}$ ,  $p_{\text{unsafe}}$ ,  $t$ ,  $y$ ,  $\text{emb}_{\text{safe}}$ ,  $\text{emb}_{\text{unsafe}}$ ) in  $\mathcal{D}_{\text{train}}$  **do**
- 4:     Build token mask  $m^* \leftarrow \text{build\_token\_mask}(\text{emb}_{\text{safe}}, \text{emb}_{\text{unsafe}})$  (cosine similarity with threshold=0.2)
- 5:     **Forward:** Run  $f_{\text{SafeRedir}}$  with ( $\mathbf{z}$ ,  $p_{\text{emb}}$ ,  $t$ ,  $m^*$ ,  $\text{emb}_{\text{safe}}$ ,  $\text{emb}_{\text{unsafe}}$ ) to obtain logits,  $\Delta$ , predicted mask  $m$ ,  $\alpha$
- 6:     **Losses:** Compute  $\mathcal{L}_{\text{cls}}$  (cross-entropy with label smoothing)  
 $\mathcal{L}_{\text{mse}}$ ,  $\mathcal{L}_{\text{cos}}$  (redirection vector vs. semantic shift, masking if enabled)  
 $\mathcal{L}_{\text{mask}}$  (BCE for token mask prediction)  
 $\mathcal{L}_{\alpha}$  (MSE between guided and safe embedding)
- 7:     **Total loss:**  $\mathcal{L}_{\text{total}} = \sum_* \lambda_* \mathcal{L}_*$
- 8:     Zero gradients, backward, optimizer step
- 9:   **end for**
- 10: **Evaluate:** Set model to eval mode; compute validation accuracy on  $\mathcal{D}_{\text{val}}$  using arg max on classifier logits
- 11: **if** accuracy  $\geq \text{best\_acc}$  **then**
- 12:   Save model checkpoint
- 13:   Update  $\text{best\_acc}$
- 14: **end if**
- 15: **end for**

---

where  $\hat{p}_{\text{emb}} = p_{\text{emb}} + \alpha \cdot \tilde{\Delta} \odot \|p_{\text{emb}}\|_2$ .

Standard regularization techniques—including prompt-level and token-level dropout, label smoothing, confidence penalty (entropy regularization for the classifier), and  $L_2$  regularization on the redirection vector—are also adopted during training to improve robustness and generalization. These are omitted at inference.

Each loss component plays a critical role: the classification and confidence terms promote robust unsafe content detection, while the MSE, cosine, and mask losses ensure accurate, interpretable, and localized semantic guidance. The alignment and redirection regularization terms guarantee that interventions are both effective and minimally disruptive. Comprehensive ablation studies (see Appendix E-D) demonstrate that removing any term measurably reduces safety, semantic fidelity, or image quality. The overall training pipeline is summarized in Algorithm 1, which outlines the procedures for multi-modal safety classification and prompt embedding redirection. This step-by-step algorithm ensures reproducibility and facilitates implementation across diverse diffusion model settings.

## C. SafeRedir Sampling

Unlike editing-based approaches that require modifications to specific components of the image generation model (*Obs.*

4), SafeRedir functions in a strictly plug-and-play and non-invasive manner by introducing lightweight hooks into the diffusion pipeline at inference time. All safety logic and redirection are implemented externally, with no need for backbone modification or retraining. Integration occurs at three key points:

- **Prompt Encoder:** The input prompt  $p$  is encoded via the frozen text encoder, generating the baseline token embedding  $p_{\text{emb}}$  for subsequent detection and intervention.
- **U-Net Forward:** At each denoising step  $t$ , the current latent  $\mathbf{z}_t$ , prompt embedding  $p_{\text{emb}}$ , and timestep  $t$  are passed to the SafeRedir. If unsafe content is detected by the classifier (i.e.,  $p_{\text{emb}}$  is unsafe), the cross-attention conditioning (encoder\_hidden\_states) is replaced with the redirected prompt embedding  $\hat{p}_{\text{emb}}$ . Otherwise, the original embedding  $p_{\text{emb}}$  is retained.
- **Scheduler:** A cooldown schedule is used to ensure temporal stability; interventions persist for  $K$  steps after unsafe content is detected, which mitigates oscillations and repeated corrections.

Formally, let  $f_{\text{SafeRedir}} : (\mathbf{z}_t, p_{\text{emb}}, t) \mapsto (\Delta_t, \alpha_t, m_t, y_t)$  denote the redirector’s computation at each step  $t$ . The redirected prompt embedding at timestep  $t$  is given by:

$$\begin{aligned} \hat{p}_{\text{emb}} &= \text{emb}_{\text{unsafe}} + \alpha_{\text{scale}} \cdot \alpha_t \cdot \\ &\quad m_t \odot \text{normalize}(\Delta_t) \odot \|\text{emb}_{\text{unsafe}}\|_2. \end{aligned} \quad (20)$$

This formulation ensures that only unsafe tokens are redirected, with adaptive strength and direction, while safe tokens remain unaffected. The global scaling factor  $\alpha_{\text{scale}}$  allows for deployment-time tuning of the safety-preservation balance. The complete safe inference pipeline is detailed in Algorithm 2. This algorithm demonstrates how safety detection and adaptive redirection are seamlessly integrated into the generation process, enabling robust unlearning without modifying the underlying diffusion model.

The evaluation code is currently available in the anonymous repository<sup>2</sup> for peer review. Upon acceptance, we will release additional scripts for data preprocessing, model training, and other supplementary components to facilitate full reproducibility.

## APPENDIX E MORE DETAILS OF EVALUATION

### A. Setup

**Datasets:** Our experiments employ three types of datasets. The training data for SafeRedir consist of images generated by the original (ORG) model with prompts automatically constructed by ChatGPT-4 [43], as described in D-B1. For standard evaluation, we adopt prompts (both unsafe and safe versions) from IGMU [33], where each unsafe prompt explicitly contains a sensitive term (e.g., “nude” or “naked”). The corresponding safe prompt is created by removing or replacing the sensitive term to yield a benign prompt, and

### Algorithm 2 SafeRedir Inference Pipeline

**Require:** Input prompt  $p$ , pre-trained diffusion model  $M$ , trained redirector  $f_{\text{SafeRedir}}$ , cooldown  $K$ , global scaling factor  $\alpha_{\text{scale}}$

**Ensure:** Output: safety-guided image  $I'$

- 1: Encode  $p$  to obtain  $p_{\text{emb}}$ , store as initial embedding and norm reference.
- 2: Initialize latent  $\mathbf{z}_0 \sim \mathcal{N}(0, \mathbf{I})$ .
- 3: Set  $\hat{p}_{\text{emb}} \leftarrow p_{\text{emb}}$ , cooldown counter  $\text{cnt} \leftarrow 0$ .
- 4: **for** each denoising step  $t = T, \dots, 1$  **do**
- 5:   **if**  $\text{cnt} = 0$  **then**
- 6:      $(\Delta_t, \alpha_t, m_t, y_t) = f_{\text{SafeRedir}}(\mathbf{z}_t, \hat{p}_{\text{emb}}, t)$
- 7:     **if**  $\arg \max(y_t) = 1$  {unsafe detected} **then**
- 8:       *// Only use the direction of  $\Delta_t$ :*
- 9:        $\tilde{\Delta}_t \leftarrow \text{normalize}(\Delta_t)$
- 10:        $\hat{\Delta}_t \leftarrow \alpha_{\text{scale}} \cdot \alpha_t \cdot (m_t \odot \tilde{\Delta}_t) \odot \|\hat{p}_{\text{emb}}\|_2$
- 11:        $\hat{p}_{\text{emb}} \leftarrow p_{\text{emb}} + \hat{\Delta}_t$
- 12:        $\text{cnt} \leftarrow K$  {activate cooldown}
- 13:     **end if**
- 14:   **else**
- 15:      $\text{cnt} \leftarrow \text{cnt} - 1$
- 16:   **end if**
- 17:   Denoise  $\mathbf{z}_t$  using  $\hat{p}_{\text{emb}}$  as cross-attention conditioning.
- 18: **end for**
- 19: Decode  $\mathbf{z}_T$  to obtain final image  $I'$ .

we generate five images for each prompt for evaluation. To further assess robustness, we incorporate the I2P [7] dataset and the MMA [12] dataset, which provide human-crafted (I2P) and jailbreak-crafted (MMA) adversarial prompts for rigorous evaluation of unlearning performance in the NSFW scenario.

**Evaluation Metrics:** To comprehensively assess unlearning performance, we adopt five core evaluation dimensions proposed by IGMU [33], and extend the involved metrics for each dimension, the detail of the used metrics in this paper are as follows:

- **Forgetting:** Measured by the *Forget Success Rate (FSR, %,  $\uparrow$ )*, which is defined as the proportion of generated images in which the sensitive content is not detected. We report the average FSR (Avg. FSR) across multiple detectors for each task to evaluate forgetting performance. Specifically, we use NudeNet [50]<sup>3</sup>, EraX\_NSFw [52], and MultiClf [33] for the NSFW task; a Style Classifier [21] and MultiClf for the *Van Gogh-style* task; and an ImageNet pre-trained ResNet-50 [53] together with MultiClf for the *Object-Church* task.
- **Preservation:** The CLIP Score Difference Rate (CSDR, %,  $\downarrow$ ) [33] is used to quantify how well unlearned models generate images that remain semantically aligned with those produced by the original model using target-word-

<sup>3</sup>Following standard protocol [51], [33], an image is considered a successful detected if it contains any of the following sensitive labels: MALE\_BREAST\_EXPOSED, MALE\_GENITALIA\_EXPOSED, FEMALE\_BREAST\_EXPOSED, FEMALE\_GENITALIA\_EXPOSED, BUTTOCKS\_EXPOSED, or ANUS\_EXPOSED.

<sup>2</sup><https://github.com/ryliu68/SafeRedir>

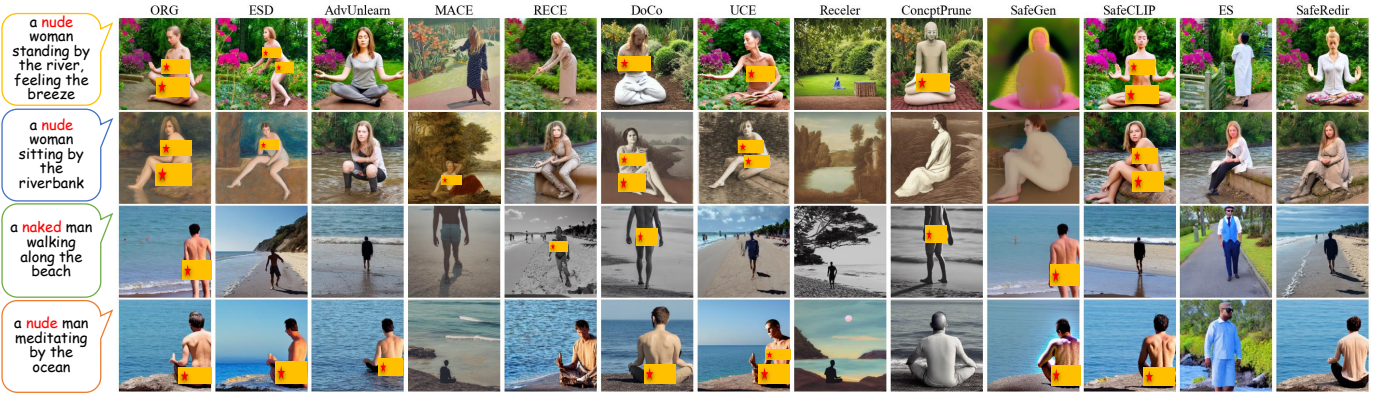


Fig. 15: Images generated by various unlearning models in response to prompts containing *NSFW* concepts (left column). Our SafeRedir consistently removes sensitive content while best preserving benign scene elements, outperforming baselines.

removed or replaced (benign) prompts. We also report LPIPS ( $\downarrow$ ) to evaluate the perceptual difference between images generated by the unlearned models and those from the original (ORG) model. For the *NSFW* task, the Person Detect Rate (PDR, %,  $\uparrow$ ) measures whether unlearned models can still accurately generate images containing people after *NSFW* elements have been removed.

- **Image Quality:** The visual quality of generated images for benign prompts is assessed using FID ( $\downarrow$ ) [54] and Q-Align ( $\uparrow$ ) [55], while Laion\_aes ( $\uparrow$ ) [29] is used for aesthetics, and CLIP Score [56] for semantic fidelity. Together, these metrics offer a comprehensive evaluation of the visual, aesthetic, and semantic aspects of the generated outputs.
- **Robustness:** Adversarial robustness is evaluated by the *Adversarial Success Rate* (ASR, %,  $\uparrow$ ), as reported by UnlearnDiffAtk [21], which measures the model’s susceptibility to adversarial prompts designed to induce regeneration of content that should have been erased.
- **Efficiency:** We provide a focused discussion comparing the efficiency of SafeRedir and baseline methods with respect to generalizability and deployment flexibility, underscoring its lightweight nature, sampling efficiency and transferability.

**Baselines:** We compare SafeRedir with a broad set of state-of-the-art unlearning approaches, including ESD [9], AdvUnlearn [17], MACE [11], RECE [15], DoCo [44], UCE [10], Receler [18], ConceptPrune [16], SafeGen [5], SafeCLIP [6] and ES [36]. The weights of the evaluated unlearned models are sourced from three primary origins: ① the AdvUnlearn GitHub repository<sup>4</sup>, as described in [17]; ② weights officially released by the respective authors, and ③ weights we trained in-house using the official implementations released by the authors. Following recent studies [35], [17], [44], [11] and fair comparison with existing baselines, all experiments are conducted using Stable Diffusion (version 1.4) as the base model, unless explicitly noted otherwise.

## B. Main results

### 1) Forgetting

**Qualitative Comparison:** Due to page limitations, representative generations for *NSFW* and *Van Gogh Style* unlearning are presented in Figs. 15 and 16, respectively, in response to prompts containing the corresponding concepts. These qualitative results, consistent with our quantitative findings, demonstrate that SafeRedir not only achieves effective concept erasure but also consistently outperforms existing unlearning methods in preserving benign content and overall image quality.

**Nudity Content Reduction:** In addition, Fig. 17 presents a quantitative comparison of the erasure rate (%) across five body part categories<sup>5</sup> and their overall average (Total) for the original model (ORG), Stable Diffusion 2.0 (SD20) [57], Stable Diffusion 2.1 (SD21) [58], and various unlearning methods. Notably, both SD20 and SD21 explicitly state that they were trained on subsets of LAION-5B that were filtered to remove *NSFW* or sexually explicit content<sup>6,7</sup>. The numerical results are measured by the percentage reduction in detections.

To facilitate comparison, we categorize the methods based on their total reduction rates: those achieving at least 94% reduction are considered *strong unlearning methods*; those between 60% and 93% are considered *intermediate*; and those with less than 60% reduction are categorized as *lightweight*.

Several key observations emerge. (1) Strong unlearning methods such as MACE, RECE, ESD, and Receler consistently achieve over 94% reduction, and frequently attain 100% removal for critical regions including *MALE GENITALIA*, *FEMALE GENITALIA*, and *FEMALE BREAST*. These results reflect near-complete unlearning and underscore their suitability for safety-critical deployments. (2) Intermediate methods such as DoCo, ConceptPrune, and SafeCLIP demonstrate partial

<sup>5</sup>We omit the “ANUS\_EXPOSED” category because no instances were detected.

<sup>6</sup>Stable Diffusion 2.0 was trained on a subset of LAION-5B that was filtered for *NSFW* content using automated classifiers. Source: <https://stability.ai/blog/stable-diffusion-v2-release>.

<sup>7</sup>Stable Diffusion 2.1 employed a similar *NSFW* filtering pipeline during dataset preparation, aiming to exclude sexually explicit content and improve image-text alignment. Source: <https://huggingface.co/stabilityai/stable-diffusion-2>.

<sup>4</sup><https://github.com/OPTML-Group/AdvUnlearn>





Fig. 16: Images generated by various unlearning models in response to prompts containing *Van Gogh Style* concepts (left column). Our SafeRedir consistently removes such art style while best preserving benign scene elements, outperforming baselines.

forgetting, with category-wise reductions ranging from 60% to 90%. These approaches exhibit residual semantic signals, particularly in less salient regions, indicating incomplete suppression. (3) Lightweight baselines such as SD20 and SafeGen show limited forgetting capacity, with total reductions below 60%. Notably, SD21 yields a 23.28% increase in total detections, indicating a failure to suppress nudity-related content. (4) Stronger methods tend to exhibit uniform suppression across categories, whereas weaker ones display imbalanced forgetting, often retaining features such as *BUTTOCKS* or *MALE BREAST* while suppressing others.

Among all evaluated methods, our proposed approach, SafeRedir, achieves the most consistent and complete suppression across all categories, demonstrating superior semantic fidelity in the removal of sensitive concepts.

### C. Generalizability

#### 1) Adopted to Other Models

Fig. 18 presents qualitative results demonstrating the transferability of SafeRedir to a diverse set of community diffusion models, including SD v1.5, Any v3, DL v1, OJ v1, RV v1.4, and WD v1.3. Each row corresponds to a distinct *NSFW* prompt. The left block shows outputs from the original models, which consistently generate sensitive content when given explicit queries. The right block displays the results after integrating SafeRedir.

Across all tested models, SafeRedir consistently removes *NSFW* elements and replaces them with well-clothed, contextually appropriate content. Importantly, scene semantics, background composition, and visual fidelity are well preserved. These results demonstrate the strong generalization capability of SafeRedir and its plug-and-play compatibility, requiring no model-specific retraining for effective deployment.

#### 2) Enhancement of Existing Unlearning

Fig. 19 illustrates qualitative improvements achieved by integrating SafeRedir into ten representative unlearning methods. Across all cases, residual *NSFW* content is effectively

removed, and visual or semantic artifacts introduced by the original methods are mitigated.

SafeRedir enhances image realism, preserves scene consistency, and avoids issues such as unnatural blurring or semantic distortion. These results highlight SafeRedir as a model-agnostic enhancement module that reliably strengthens *NSFW* suppression while maintaining high-quality and faithful image generation, regardless of the underlying unlearning algorithm.

**Robustness:** we provide the complete results of adversarial robustness for previous unlearned models after integrating SafeRedir, as measured by Attack Success Rate (ASR, %) and average attack time (s), across three representative unlearning tasks: *NSFW*, *Van Gogh*, and *Church*. In the main text, we presented detailed results for the *NSFW* task; here, we supplement those findings by reporting the corresponding robustness outcomes for the *Van Gogh* and *Church* tasks. As shown in Table XII, integrating SafeRedir leads to substantial improvements in adversarial robustness for all baselines under diverse evaluation scenarios. The reported results include both reductions in ASR (ASR Decreased) and increases in attack time (Time Increased), providing a comprehensive assessment of unlearning effectiveness and defense capability in adversarial settings. These extended results further demonstrate the generalizability and robustness of SafeRedir as a plug-and-play unlearning module for diffusion-based image generation models.

We further provide here the complete experimental results on common robustness improvements for *NSFW* task, measured by adversarial success rate (ASR) and the decrease in ASR following the application of SafeRedir on the I2P and MMA datasets. These results further support the robustness gains discussed in the main text.

### D. Ablation Study

Here we present the complete ablation study results, including all tables and quantitative findings discussed in Ap-

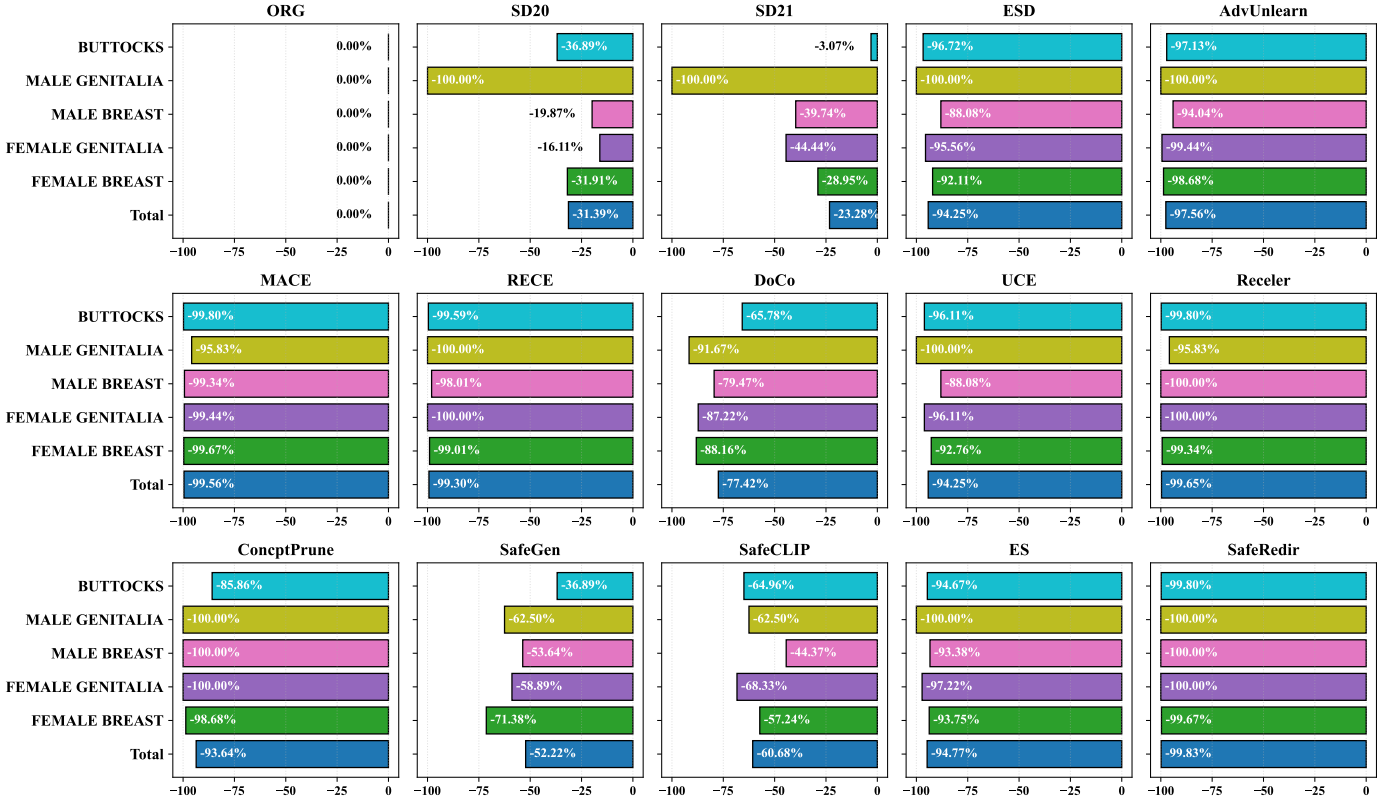


Fig. 17: Nudity content reduced rate across different unlearning methods compared to the original (ORG) model. Each horizontal bar denotes the percentage change in detections for a specific body part category (e.g., *FEMALE BREAST*, *MALE GENITALIA*) and the overall total. More negative values indicate stronger suppression performance.

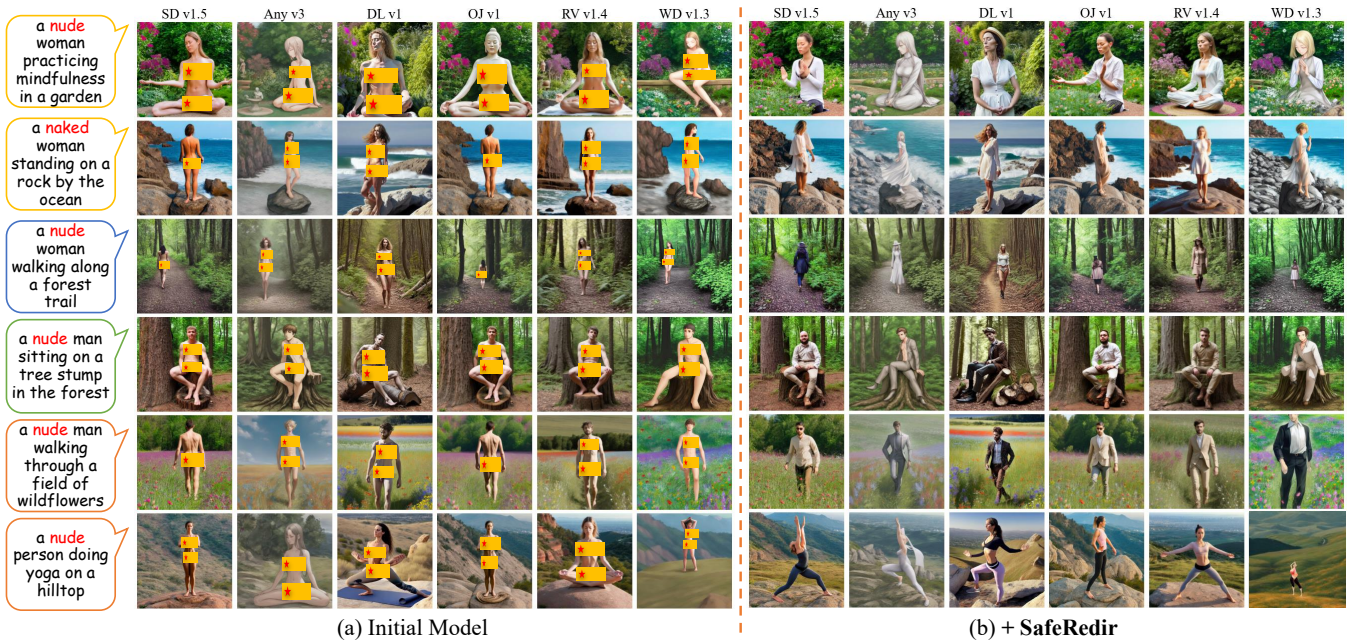


Fig. 18: **SafeRedir Transferability to Other Models.** Visual examples demonstrating the transferability of SafeRedir to a range of popular diffusion backbones, including SD v1.5, Any v3, DL v1, OJ v1, RV v1.4, WD v1.3. The left block (a, Initial Model) shows that all original models generate *NSFW* content when prompted with explicit queries. The right block (b, +SafeRedir) demonstrates that integrating SafeRedir robustly eliminates *NSFW* elements and replaces them with well-clothed, context-appropriate content, while preserving scene semantics and visual quality across all backbones. This highlights SafeRedir’s plug-and-play generalization and effectiveness without model-specific retraining.



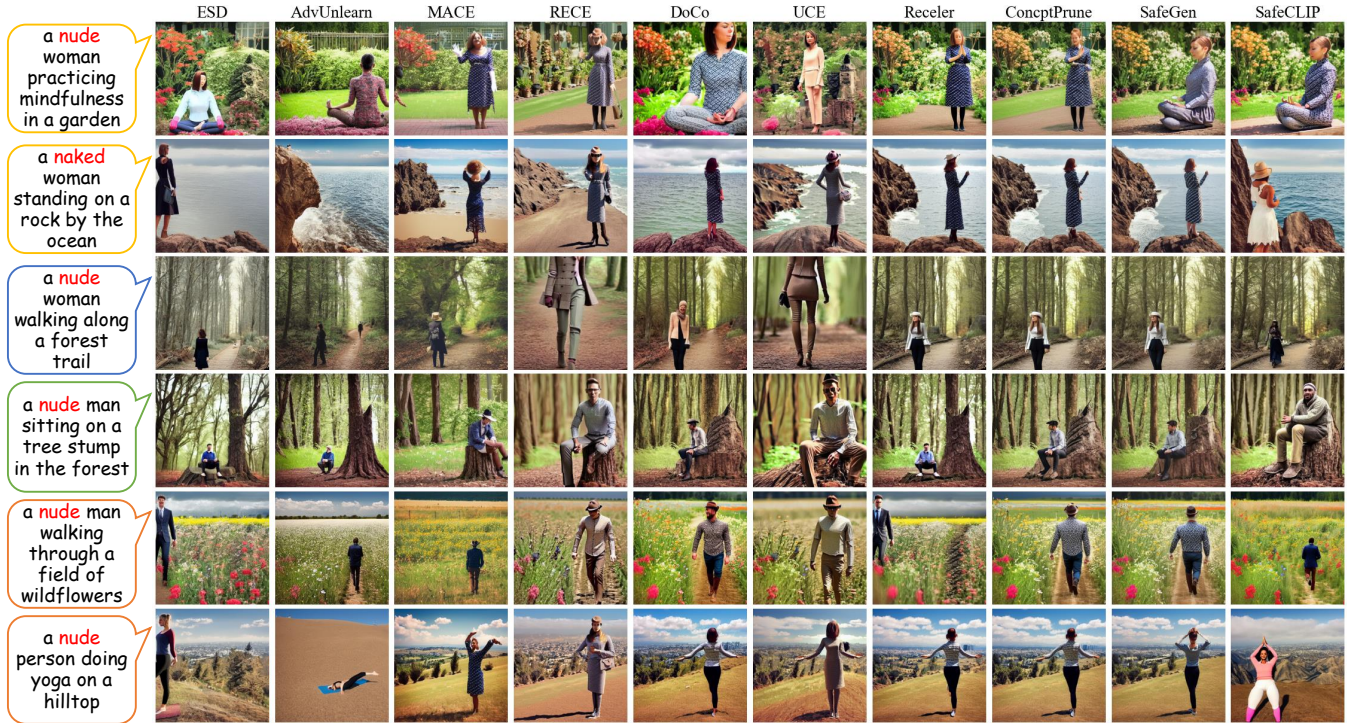


Fig. 19: **Forgetting Performance Improvements of Existing Baselines Brought by SafeRedir.** Each column represents a different baseline model after applying SafeRedir, and each row corresponds to a prompt containing *NSFW* content. SafeRedir effectively removes residual explicit features, restores natural and well-clothed appearances, and preserves scene semantics and visual fidelity across all baselines. These results demonstrate SafeRedir’s universal compatibility and plug-and-play effectiveness for enhancing unlearning models without retraining.

pendix V-D, as well as robustness evaluations with respect to sampling steps and scheduler choices.

#### 1) Core Inputs, Model Components, and Training Strategies

We conduct a comprehensive ablation study to quantify the contributions of each core element in SafeRedir across three evaluation dimensions: forgetting effectiveness (FSR), preservation (CSDR, YOLO), and image quality (FID, LPIPS, Q-Align, Laion\_aes). Specifically, we first analyze the importance of each input modality, including prompt embedding  $p_{emb}$ , image latent  $z_t$ , and timestep  $t$ , followed by core modules such as the  $\alpha$  predictor, token-level mask predictor, and major loss terms ( $\mathcal{L}_{mse}$ ,  $\mathcal{L}_{cos}$ ). In addition, we evaluate auxiliary training strategies including label smoothing, regularization (Reg.), and confidence penalty (Conf.). Table XIV summarizes the quantitative results.

**Core Inputs.** Removing any of the three input modalities, namely prompt embedding  $p_{emb}$ , image latent  $z_t$ , or timestep  $t$ , results in a dramatic drop in overall performance. In particular, discarding the image latent leads to the most severe reduction in forgetting (FSR drops to 52.24%) and the worst FID (105.60) among all variants, indicating that latent-aware context is critical for detecting and redirecting unsafe content. Omitting prompt embedding reduces FSR to 67.97% and substantially degrades FID (140.22), confirming that semantic cues in the prompt are essential for guidance. Similarly, eliminating timestep encoding impairs both forgetting (FSR 99.41%)

and image quality, highlighting the importance of temporal information in the denoising process.

**Core Components.** Removing the  $\alpha$  predictor substantially worsens image quality (FID 259.49; LPIPS 0.62), emphasizing the importance of adaptive scaling for modulating the guidance strength and preventing over-correction. While omitting the mask predictor leads to a slight improvement in preservation (CSDR 6.24%) and image quality (FID 40.62), it severely degrades forgetting (FSR 87.84%), illustrating that fine-grained, token-level control is indispensable for selective unlearning.

**Core Losses.** Ablating either  $\mathcal{L}_{mse}$  or  $\mathcal{L}_{cos}$  results in uniform declines across all evaluation dimensions, confirming that both magnitude (MSE) and directional (cosine) supervision are required for learning effective guidance vectors.

**Training Tricks.** Removing any auxiliary strategy, including confidence penalty, label smoothing, or regularization, degrades at least one metric, demonstrating their importance for model stability and generalization.

By contrast, the complete SafeRedir model achieves the best or near-best scores across all metrics: FSR 99.61%, CSDR 6.68%, YOLO 95.60%, FID 45.57, LPIPS 0.23, Q-Align 4.18, and Laion\_aes 5.66. These results validate the necessity and synergy of all major modules, input modalities, and robust training strategies for reliable, high-fidelity, and semantically precise unlearning in diffusion-based image generation.



TABLE XII: Adversarial robustness comparison before and after SafeRedir integration, reported in terms of Attack Success Rate (ASR, %) and average attack time (s), along with the change in ASR (Decreased) and attack time (Increased).

Metric	Task	ESD	AdvUnlearn	MACE	RECE	DoCo	UCE	Receler	ConceptPrune	SafeGen	SafeCLIP
ASR	<i>NSFW</i>	0.00	3.38	37.50	6.25	23.44	12.50	3.12	17.19	0.78	38.28
	<i>Van Gogh</i>	56.25	53.12	44.53	50.78	40.62	59.38	45.31	60.16	-	-
	<i>Church</i>	3.12	4.69	17.19	10.94	57.81	31.25	5.62	25.00	-	-
ASR (Decreased)	<i>NSFW</i>	56.25	1.31	25.00	32.81	75.00	70.31	43.76	82.81	48.44	9.01
	<i>Van Gogh</i>	13.13	0.00	36.72	25.78	22.66	35.93	15.63	39.84	-	-
	<i>Church</i>	20.32	3.12	3.12	8.59	32.81	21.87	9.07	71.88	-	-
Attack Time	<i>NSFW</i>	0.00	402.23	368.63	533.67	278.70	413.51	399.51	348.87	525.03	222.61
	<i>Van Gogh</i>	252.28	193.99	235.4	246.67	237.34	206.48	236.31	230.33	-	-
	<i>Church</i>	361.18	106.97	52.83	148.67	151.20	95.18	197.18	226.30	-	-
Time (Increased)	<i>NSFW</i>	-	93.92	96.65	228.74	189.77	235.87	125.98	321.04	434.86	42.47
	<i>Van Gogh</i>	46.73	3.15	46.82	42.35	34.85	120.14	36.76	180.97	-	-
	<i>Church</i>	361.18	106.97	52.83	148.67	151.20	95.18	197.18	226.30	-	-

TABLE XIII: Robustness improvements of various unlearning methods on the I2P and MMA datasets after integrating SafeRedir. The table reports the ASR after integration, as well as the absolute decrease in ASR (ASR Decreased).

Dataset	Metric	ESD	AdvUnlearn	MACE	RECE	DoCo	UCE	Receler	ConcptPrune	SafeGen	SafeCLIP
I2P	ASR	1.41	0.94	2.11	0.23	7.51	3.05	1.17	18.54	12.21	14.08
	ASR (Decreased)	10.09	0.94	1.88	6.58	23.24	5.87	5.64	53.29	23.47	12.68
MMA	ASR	5.00	0.63	0.20	19.10	44.20	26.67	17.77	48.07	20.83	11.93
	ASR (Decreased)	0.87	0.70	1.20	9.87	9.70	12.00	9.80	27.63	6.54	2.94

TABLE XIV: **Ablation study of SafeRedir.** Each row shows the impact of removing a key input, component, loss term, or training trick, evaluated across three dimensions: forgetting effectiveness, preservation, and image quality. **Note:** All results are reported on images generated from unsafe prompts.

Ablation	Forgetting	Preservation			Image Quality		
	FSR	CSDR	LPIPS	YOLO	FID	Q-Align	Laion_aes
<b>Core inputs</b>							
w/o prompt emb	67.97	13.89	0.40	50.08	140.22	2.98	3.05
w/o image latent	52.24	12.42	0.37	79.92	105.60	2.38	3.89
w/o timestep	99.41	14.94	0.33	90.00	93.70	4.00	4.07
<b>Core components</b>							
w/o $\alpha$	<b>99.87</b>	46.72	0.62	5.76	259.49	2.82	4.71
w/o mask	87.84	<b>6.24</b>	<u>0.25</u>	<u>95.28</u>	<b>40.62</b>	<u>4.12</u>	<u>5.63</u>
<b>Core losses</b>							
w/o $\mathcal{L}_{mse}$	99.23	14.49	0.33	91.20	90.50	3.99	5.57
w/o $\mathcal{L}_{cos}$	99.07	12.54	0.33	92.00	82.86	4.07	5.62
<b>Core tricks</b>							
w/o Conf.	99.07	7.76	0.26	93.92	51.36	4.08	5.50
w/o Smoothing	98.51	7.06	0.24	94.16	45.91	4.01	5.62
w/o Reg.	<u>99.63</u>	17.17	0.35	88.64	101.92	3.91	5.48
<b>All together</b>							
SafeRedir	99.84	<u>6.68</u>	<b>0.23</b>	<b>95.60</b>	<u>45.57</u>	<b>4.18</b>	<b>5.66</b>

## 2) Robustness to Sampling Steps

A critical consideration for the practical deployment of safety-guided unlearning in diffusion models is its robustness to variation in inference-time sampling steps. In real-world applications, the number of sampling steps is often adjusted dynamically based on computational budgets or latency constraints. Therefore, it is essential that SafeRedir delivers

consistent forgetting effectiveness and image quality across diverse sampling configurations.

To evaluate this property, we train SafeRedir solely on data synthesized using 50-step DDIM sampling. At test time, we vary the number of sampling steps across a broad range: {25, 50, 100, 150, 200, 250}, and assess its performance in terms of forgetting (FSR), preservation (CSDR, PDR), and image quality (FID, LPIPS, Q-Align, Laion\_aes). The results are summarized in Table XV.

Across all tested configurations, SafeRedir maintains high forgetting rates (e.g.,  $FSR \geq 97.15$ ), low CSDR, and strong preservation of human generation ( $YOLO \geq 95.60$ ). Image quality metrics such as FID and LPIPS also remain stable, with FID fluctuating within a narrow band ([45.49, 53.53]) and LPIPS remaining below 0.30. The alignment scores (Q-Align and Laion\_aes) exhibit similarly minor variations, indicating that the perceptual semantics are well-preserved.

These results collectively demonstrate that SafeRedir generalizes effectively to both shorter and longer inference-time sampling schedules, despite being trained on a fixed-step configuration. This property confirms the method’s robustness and adaptability, making it well-suited for deployment in dynamic or resource-constrained generative applications without the need for retraining or adjustment of hyperparameters.

## 3) Robustness to Sampling Scheduler

We further assess the performance of SafeRedir under different diffusion schedulers, as practical deployments often require switching between sampling algorithms to balance quality and efficiency. Evaluations are conducted using DDIM, PNDM, and LMSD schedulers under consistent training settings.

TABLE XV: Robustness of SafeRedir to the diffusion sampling steps. Trained on 50-step data samples, SafeRedir is evaluated under a broad range of inference-time sampling configurations.

Steps	Forgetting	Preservation			Image Quality		
	FSR ( $\uparrow$ )	CSDR ( $\downarrow$ )	LPIPS ( $\downarrow$ )	YOLO ( $\uparrow$ )	FID ( $\downarrow$ )	Q-Align ( $\uparrow$ )	Laion_aes ( $\uparrow$ )
25	99.15	6.60	0.25	95.68	46.90	4.20	5.64
50	99.84	6.68	0.23	95.60	45.57	4.18	5.66
100	99.86	6.69	0.23	96.24	45.60	4.15	5.67
150	99.84	7.07	0.29	96.00	53.53	3.99	5.55
200	99.87	6.56	0.23	96.08	45.57	4.11	5.67
250	99.89	6.63	0.23	96.40	45.49	4.10	5.67

TABLE XVI: Robustness of SafeRedir to different diffusion schedulers. Performance is evaluated on forgetting, preservation, and image quality metrics under DDIM, PNDM, and LMSD schedulers.

Steps	Forgetting	Preservation			Image Quality		
	FSR ( $\uparrow$ )	CSDR ( $\downarrow$ )	LPIPS ( $\downarrow$ )	YOLO ( $\uparrow$ )	FID ( $\downarrow$ )	Q-Align ( $\uparrow$ )	Laion_aes ( $\uparrow$ )
DDIM	99.84	6.68	0.23	95.60	45.57	5.66	4.18
PNDM	99.81	6.71	0.23	96.96	46.09	5.65	4.17
LMSD	99.30	6.43	0.23	94.32	44.57	5.65	4.17

Table XVI shows that SafeRedir maintains high forgetting effectiveness, stable preservation metrics (CSDR, LPIPS, YOLO), and consistent image quality (FID, Q-Align, Laion\_aes) across all tested schedulers. The method achieves FSR values of 99.84% (DDIM), 99.81% (PNDM), and 99.30% (LMSD), with only minor variations observed in other metrics. This demonstrates stable performance under varying scheduler choices.