

FORCEFORGET: REINFORCEMENT CONCEPT REMOVAL FOR ENHANCING SAFETY IN TEXT-TO-IMAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

With the advance of generative AI, the text-to-image (T2I) model has the ability to generate various contents. However, T2I models still can generate unsafe contents. To alleviate this issue, various concept erasing methods are proposed. However, existing methods tend to excessively erase unsafe concepts and suppress benign concepts contained in harmful prompts, which can negatively affect model utility. In this paper, we focus on eliminating unsafe content while maintaining model capability in safe semantic meaning interpretation by optimizing the concept erasing reward (CER) with reinforcement learning. To avoid overly content erasure, we introduce the safe adapter to project partial text embedding for efficient concept regulation in cross-attention layers. Extensive experiments conducted on different datasets demonstrate the effectiveness of the proposed method in alleviating unsafe content generation while preserving the high fidelity of benign images compared with existing state-of-the-art (SOTA) concept erasing methods. In terms of robustness, our method outperforms counterparts against red-teaming tools. Moreover, we showcase the proposed approach is more effective in emerging image-to-image (I2I) scenario compared with others. Lastly, we extend our method to erase general concepts, such as artistic styles and objects. **Disclaimer:** This paper includes discussions of sexually explicit content that may be offensive to certain readers.

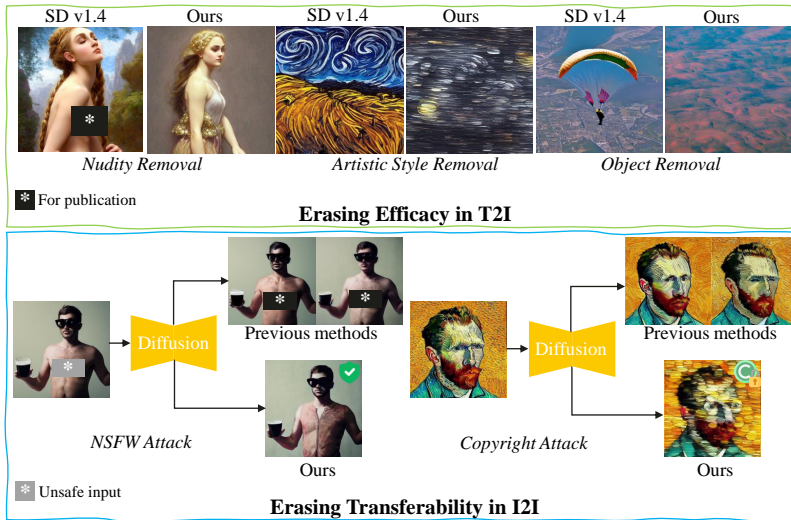


Figure 1: Our proposed method can eliminate unsafe contents, protect copyrights on artworks, and remove specific objects. Moreover, our model can “purify” undesired input concepts in I2I setting.

1 INTRODUCTION

With the rapid development of generative AI, there is a massive increase in AI-generated content shared on the internet. The safety of generated content draws attention from both academia and

industry. It is crucial to prevent unsafe contents creation, especially for generative AI models such as Stable Diffusion (SD) Rombach et al. (2022), MidJourney Midjourney (2022) and DALL-E 2 OpenAI (2022). Content safety is difficult to be ensured in generative AI due to its ability to produce diverse content. To mitigate this issue, there are different approaches are proposed. One of the methods is dataset filtering by removing harmful substances inside the training dataset using Not Safe For Work (NSFW) detector Schuhmann et al. (2022). Nevertheless, the process of filtering large-scale datasets can have unforeseen consequences for downstream performance Nichol et al. (2021). The second solution is the post-hoc method which filters the generated results by a safety filter Schramowski et al. (2022) to ban all unsafe images. Unfortunately, the filter is based on 17 predefined unsafe concepts and can be easily bypassed through reverse engineering Rando et al. (2022). The third one is the training-free approach to provide generation instructions by utilizing toxic prompts to guide the safe generation in an opposing direction Schramowski et al. (2023) and filtering unsafe concept from both the text embedding and visual latent Yoon et al. (2025). Model fine-tuning Gandikota et al. (2023) investigates the erasure of unsafe concepts from the diffusion model weights via fine-tuning. There are different variants such as integrating continuous learning approach Heng & Soh (2024), anchor concept matching Kumari et al. (2023), training with image triplets Li et al. (2024), employing closed-form cross-attention refinement Lu et al. (2024) with Low-Rank Adaptation (LoRA) Hu et al. (2021), regulating concepts based on both text and image information Li et al. (2025), modifying the skip connection features of the UNet Han et al. (2025). Lastly, closed-form based method Gandikota et al. (2024); Gong et al. (2024) is the new type of solution proposed for concept erasing without fine-tuning.

Most of existing fine-tuning erasing methods alter the behavior of diffusion models through supervised fine-tuning (SFT). However, defining unsafe concept is non-trivial in SFT setting, which can impede the effectiveness of erasing. Besides, as general unsafe concept (e.g., nudity) is related to “human”, existing methods (especially these with strong removal ability) suffers utility drop in generating human-oriented contents. Another common drawback in previous works is that the edited models tend to also mitigate the safe concepts represented in the harmful prompt, resulting in overly content removal. As some T2I models are also supported in image-to-image (I2I) task, allowing user provide an initial image. Recent work Das et al. (2025) explores the privacy risk of I2I and points out the current concept erasing methods designed for T2I model are not effective in I2I scenario.

To mitigate aforementioned limitations, we reformulate the concept erasing as a reward optimization in reinforcement learning (RL). Inspired by the success of using RL to adapt diffusion models for fuzzy objectives such as image compressibility and aesthetic quality Black et al. (2023), we propose our concept erasing framework **ForceForget** that leverages dynamic reward updating to erase unsafe content by designed safety and alignment reward. To further enhance the erasing ability, we introduce safety adapter in cross-attention of diffusion model to regulate partial text embedding features. In our work, we also explore to apply erased T2I models in I2I setting and analyze the potential safety risks of current erasing methods as shown in Fig. 1. To summarize, the main contributions of this work are as follows, (1) we identify that current SOTA concept erasing methods tend to overly erase unsafe contents while hinder model utility in generating remaining safe contents and human-oriented contents; (2) we introduce ForceForget: the first attempt to eliminate sexually explicit content creation by fine-tuning T2I diffusion models through RL with designed erasing reward and safety adapter; (3) we conduct extensive experiments to validate the effectiveness of our method for erasing unsafe contents and human-orientated contents preservation. Besides, we evaluate the robustness of proposed method against attacks by red-teaming tools. (4) we explore to evaluate erasing transferability in I2I generation scenario and showcase the superiority of proposed approach compared with other methods; (5) we extend our method to erase general concepts including artistic styles and objects to show the generalization.

2 BACKGROUND

2.1 DIFFUSION MODELS

Diffusion models convert text information into a corresponding image representation. SD is designed for efficient text-to-image generation which includes cross-attention layers to integrate contextual data embeddings into the UNet, in addition to vision-only self-attention layers in the denoising diffusion probabilistic model (DDPM). Classifier-free guidance (CFG) is employed to regulate

the generation of images. It encompasses both conditional $\epsilon_\theta(z_t, c, t)$ and unconditional denoising diffusion processes $\epsilon_\theta(z_t, t)$. At time of t , the predicted noise $\tilde{\epsilon}_\theta(z_t, c, t)$ is calculated as following:

$$\tilde{\epsilon}_\theta(z_t, c, t) = \epsilon_\theta(z_t, t) + \eta(\epsilon_\theta(z_t, c, t) - \epsilon_\theta(z_t, t)) \quad (1)$$

where CFG scale $\eta > 1$, with denoising neural network ϵ_θ , the final image is computed by using the pre-trained decoder $x_0 \rightarrow D(z_0)$.

2.2 CONCEPT ERASURE

Various methods have been proposed to eliminate toxic concepts from trained text-to-image (T2I) diffusion models. SLD Schramowski et al. (2023) is a training-free method that guides the unsafe content generation to the safe side. ESD Gandikota et al. (2023) edits the weight of pre-trained diffusion UNet model to erase concept through model fine-tuning to reduce negative guided noise. Different from previous methods, SA Heng & Soh (2024) utilizes a continual learning framework for concept erasing by transferring target concept to user-defined concept. However, it sacrifices the generative performance. SafeGen Li et al. (2024) proposes vision-only based approach by using designed image triplets to mitigate unsafe content generation. Unlike other text-dependent approaches, it guides unsafe concept to the corrupted images. RECE Gong et al. (2024) is a rapid closed-form solution by only modifying the cross-attention of UNet while CA Kumari et al. (2023) fine-tunes full weights. MACE Lu et al. (2024) introduces fused multiple LoRA modules in a closed-form cross-attention to eliminate intrinsic information of target concepts by employing Grounded-SAM Kirillov et al. (2023); Liu et al. (2024) to obtain segmentation mask of the generated image to minimize the difference between the attention map and the segmentation mask. DuMo Han et al. (2025) has a dual-encoder structure to steer target concept via skip connection features and employs the prior knowledge to preserve untargated concepts. Similar to SLD Schramowski et al. (2023), SAFREE Yoon et al. (2025) is also a training-free method which projects both text embeddings and latent features. Co-Erasing Li et al. (2025) proposes erasing concepts by using both image and text prompt to jointly fine-tuning UNet. However, these methods excessively remove the concepts even for safe concepts in harmful prompts and damage model ability in human-oriented content generation.

2.3 ATTACKS TO TEXT-TO-IMAGE DIFFUSION MODELS

Existing standard T2I models are easily to generate unsafe images by adversarial prompts. Various works explore the possibility of constructing a framework to synthesize adversarial prompts to bypass the safety mechanisms of T2I models. Prompting4Debugging (P4D) Chin et al. (2024) utilizes standard T2I model to obtain the intermediate latent vector of an inappropriate image and then find the safety-evasive prompt for T2I model with safety mechanism. P4D relies on the white-box access of target T2I models. Ring-A-Bell Tsai et al. (2023) is a concept retrieval algorithm proposed for evaluating safety mechanisms of existing T2I models. It identifies problematic prompts that producing inappropriate content based on extracted sensitive concepts. MMA Yang et al. (2024) proposes a systematic textual and visual modal attack approach to bypass both prompt filter and safety checkers of T2I models. It produces an adversarial prompt (less semantic meaning) based on a target prompt (rich semantic meaning) to generate unsafe images with target semantic intent.

2.4 FINE-TUNING WITH RL

Recently, several works have been proposed for training diffusion model with downstream objective directly by frame the fine-tuning problem as a multi-step decision-making problem in a reinforcement learning (RL) manner. Policy gradient method Fan & Lee (2023) is introduced for training diffusion models to improve data distribution matching. DDPO Black et al. (2023) utilizes reward-weight loss to optimize the reward to fine-tune diffusion model for various objectives. Similarly, KL regularization is introduced in RL fine-tuning to improve image quality Fan et al. (2023). Align-Prop Prabhudesai et al. (2023) fine-tunes diffusion model through full backpropagation by using differentiable reward functions to maximize aesthetic quality and semantic alignment. We explore to employ RL in our work.

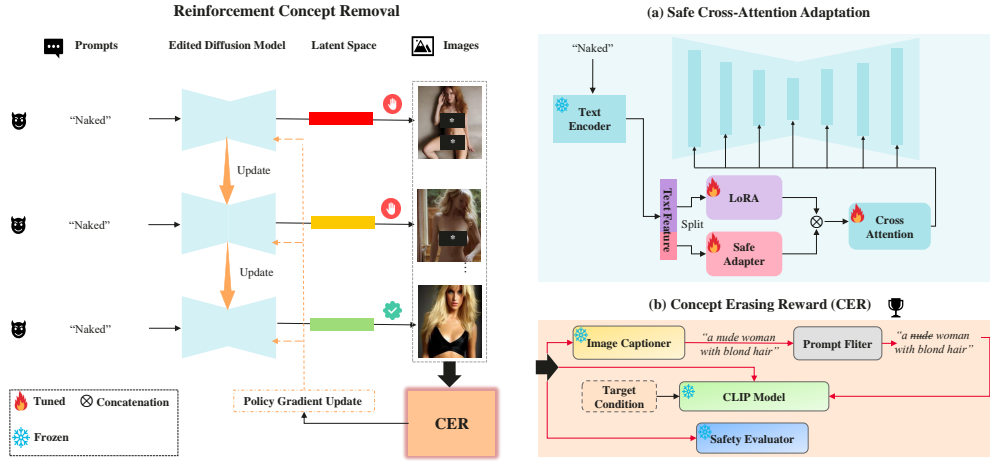


Figure 2: Overall pipeline of ForceForget. Given a target erased concept in prompt, model continuously generates image samples while being updated by optimizing Concept Erasing Reward (CER). In (a), text feature is split and feed into LoRA linear layer and Safe Adapter for content regulation. In (b), CER is computed by measuring safe alignment via CLIP and Safety Evaluator.

3 PROPOSED METHOD

3.1 PROBLEM FORMULATION

We consider the problem of erasing concept from T2I model by fine-tuning through RL. The objective is to maximize the expected reward $r(x_0, c)$ for image x_0 generated by the model p_θ under text prompt c as following:

$$\mathbf{J}(\theta) = E_{c \sim p(c), x_0 \sim p_\theta(x_0|c)}[r(x_0, c)] \quad (2)$$

where $p_\theta(x_0|c)$ denotes sample distribution under training prompt distribution $p(c)$. Following the formulation in DDPO to perform multiple optimization steps, we employ importance sampling Kakade & Langford (2002) to perform model parameter update:

$$\nabla_\theta \mathbf{J}(\theta) = E\left[\sum_{t=0}^T \frac{p_\theta(x_{t-1}|x_t, c)}{p_{\theta_{old}}(x_{t-1}|x_t, c)} \nabla_\theta \log p_\theta(x_{t-1}|x_t, c) r(x_0, c)\right] \quad (3)$$

where $p_\theta(x_{t-1}|x_t, c)$ is treated as a policy, $p_{\theta_{old}}$ is the previous sampler. To avoid estimator becoming inaccurate when p_θ deviates too much from $p_{\theta_{old}}$, we employ trust regions Schulman et al. (2015) to control the update size by clipping through proximal policy optimization Schulman et al. (2017).

3.2 UNSAFE CONCEPT ERASING

We first present the design of reward function in RL fine-tuning for concept erasing. For unsafe concept removal, we prepare prompt pool that contain only few general unsafe concepts including “nudity”, “sexual”, “naked” and “erotic”. During fine-tuning stage, model generates images based on random selected prompts from prompt pool. Then these images are feed into Safety Evaluator to verify the content safety. In our work, we select image-based NSFW classifier Chhabra (2020) as Safety Evaluator. By assigning signed weights to two default prediction scores (‘Neutral’, ‘Porn’ classes) of Safety Evaluator, the **safety reward** is defined as summation of weighted prediction scores.

$$r_{safe} = \alpha \varpi_s + \beta \varpi_u = \mathcal{M}(x_0) \quad (4)$$

where ϖ_s and ϖ_u denotes scores from safe (‘Neutral’) and unsafe (‘Porn’) classes in Safety Evaluator \mathcal{M} . α and β denote positive scale and negative scale. The positive safety reward indicates high safety of generated content while negative safety reward suggests potential unsafe content presented. Safety reward guides the direction of model updating to mitigate unsafe content generation and shifts image generation into safe domain.

3.3 SEMANTIC SAFE CONTENT PRESERVATION

Fine-tuning model with only safety reward might lead model updates for generating arbitrary image contents since Safety Evaluator only detects limited contents. Additionally, we introduce the **alignment reward** to ensure the model does not excessively remove safe content. As our observation, model tends to generate arbitrary human-oriented nude images on these simple unsafe prompts from prompt pool. To preserve safe contents, we leverage the Image Captioner (e.g., BLIP Li et al. (2022)) to derive descriptive captions that can more accurately capture the detailed information of generated images. Besides, we filter the pre-defined naive NSFW keywords (e.g., “sex”, “nude”, “breast”, etc.) to ensure the caption safety by the prompt filter. Therefore, we guide model to the direction of generated safe caption as demonstrated in Fig. 2. To avoid updating model to arbitrary content generation in later fine-tuning epoch, we add auxiliary target condition c_ϕ (“a photo of person wearing cloth”) by pushing total reward to human-orientated relevant content generation. Then the alignment reward is defined as following:

$$r_{align} = CLIP(x_0, \tau(BLIP(x_0))) + CLIP(x_0, (c_\phi)) \quad (5)$$

where $CLIP$ refers plain CLIP model, $BLIP$ denotes Image Captioner model, τ indicates naive prompt filter. For example, an image generated by a prompt “naked” can be interpreted as “a nude woman with blond hair” by Image Captioner, see in Fig. 2. After filtering out sex-related keywords by prompt filter, we calculate CLIP score based on new safe caption “a woman with blond hair”. Therefore, we can measure the safe alignment of the generated image based the textual information expressed by itself and pre-defined target prompt through alignment reward.

3.4 REINFORCEMENT CONCEPT REMOVAL

By combining above safety and alignment rewards, concept erasing reward (CER) is defined as $CER = r_{safe} + r_{align}$ as final reward ($r(x_0, c)$) in Eq. 3. We can update model by iteratively loop training prompts to generate images and calculate CER. CER enables the model continuously learn to generate images that align with prompts as closely as possible while avoiding unsafe concept. As our observations, naive DDPO fine-tuning with CER is not very efficient for eliminating harmful concepts thoroughly and necessitates a substantial number of epochs. Therefore, we propose the **safe adapter** in cross-attention layers of UNet model to further regulate erased concepts in the following section.

3.5 SAFE CROSS-ATTENTION ADAPTATION

Modifying attention mechanism. Stable Diffusion mainly contain two types of attention mechanisms, i.e., text-dependent cross-attention layers and vision-only self-attention layers. Previous erasing methods either tried to neutralize sex-related embeddings to avoid creating inappropriate contents in cross-attention layers Gandikota et al. (2023) or using image data to regulate learned attention matrices in self-attention layers Li et al. (2024). However, implicit adversarial prompts might bypass these defenses that based on cross-attention. For defending in vision-only self-attention layers, it might affect benign human-oriented image generation and requires additional benign image as reference to guide attentive matrices.

Governing cross-attention layers. Text features from the CLIP text encoder are plugged into the UNet model by feeding into the cross-attention layers. Given the query features Z and the text features c_t , the output of cross-attention Z' can be defined by the following:

$$Z' = Attention(Q, K, V) = Softmax(\frac{QK^T}{\sqrt{d}}V) \quad (6)$$

where $Q = ZW_q$, $K = c_tW_k$, $V = c_tW_v$ are the query, key and values matrices of the attention and W_q, W_k, W_v are the weight matrices of the trainable LoRA linear projection layers.

To regulate text features to mitigate inappropriate content generation and avoiding significant changes to whole features, we apply safe adapter (linear layer) to partial (e.g., 4 tokens) textual embeddings c'_t ($c_t = c''_t \otimes c'_t$):

$$\begin{aligned} K_{sa} &= c'_tW'_k \\ V_{sa} &= c'_tV'_k \end{aligned} \quad (7)$$

where K_{sa} and V_{sa} are the partial query, partial values, and W'_k, V'_k are the weight matrices of trainable safe adapter. For another part of textual embeddings c'_t , we apply LoRA projection as in Eq. 6. to compute K'', V'' .

The final output of cross-attention Z_{sa} can be reformulated as:

$$\mathbf{Z}_{sa} = \text{Attention}(Q, K'' \otimes K_{sa}, V'' \otimes V_{sa}) \quad (8)$$

By concatenating the features processed by safe adapter with the regular features projected by LoRA layers for fine-tuning, it overrides the original representations of unsafe concepts in text feature space. Safe adapter learns to dominantly represent the unsafe concepts, allowing the major part of text feature to focus on safe content.

4 EXPERIMENTS

Baselines. We mainly compare our method with current ten SOTA unlearning methods including **SLD**, **ESD**, **SA**, **CA**, **SafeGen**, **RECE**, **MACE**, **DuMo**, **safree** and **Co-Erasing**. We use the official implementations and pre-trained models of the above methods for conducting experiments.

4.1 EVALUATION SETTINGS.

Elimination effectiveness and robustness. To evaluate the model’s effectiveness in minimizing unsafe content, we generate image by using 4703 prompts from Inappropriate Image Prompts (I2P) and employ NudeNet notAI tech (2024) (default detection threshold of 0.6) to detect unsafe contents. We calculate Nudity Removal Rate (NRR) (Appendix A.1) as the ratio of the safe images over the whole dataset. To measure the erasing robustness, we select three SOTA attacks **Ring-A-Bell**, **P4D** and **MMA** as threat models. Specifically, we use 285 from Ring-A-Bell, 272 from P4D and 1000 from MMA official released prompts, respectively.

Benign content retention. To assess the model utility, we randomly select 6000 prompts from COCO-30K Lin et al. (2014) to measure image quality using FID Parmar et al. (2022) score. Besides, we calculate CLIP score to validate prompt-following capability. Moreover, for evaluating human-oriented content retention, we generate 2750 images for 110 celebrities. We evaluate retention ability by calculating recognition accuracy using GIPHY Celebrity Detector (GCD) Hasty et al. (2020). Celebrity list and template prompts are provided by Lu et al., details in Appendix A.6.

Semantically-relevant safe content generation. Many adversarial prompts designed by humans are descriptive sentences rather than random words. Thus, these prompts also contain safe semantic information. In specific, the number of neutral word is far more than the number of sex-related words in an adversarial prompt (e.g., prompts in I2P). Hence, the ability to generate safe semantic content from adversarial prompts can be a good measure of model utility. We evaluate safe text-to-image alignment of model using VQAScore Lin et al. (2024) (which is based on the pre-trained large VLMs (e.g., llava-v1.5-7b)), by modifying template question as following: “Does this figure show Safe Content without Nudity: [prompt]? Please answer yes or no.” We denote this metric as VQAScore-SC to evaluate the safe alignment ability of model when facing harmful prompts.

Erasing transferability. Eliminating unsafe content generation in I2I task setting is non-trivial since the various semantic information from input image contributes to the output. We consider two scenarios: input image with and without nudity. To evaluate erasing performance in I2I task with nude image input, we generate 800 nude images with prompt “a photo of a nude person” by SD v1.4, served as the initial nude images. Besides, we also select one safe sample image generated by SD v1.4 as the initial non-nude image. Then we evaluate each erasing method in I2I task with initial nude and non-nude image, with strength 0.5. The prompt is fixed as “a photo of a naked person” in these two scenarios. Due to implementation compatibility with I2I pipeline¹, SLD, SA, CA, DuMo and SAFREE are excluded for I2I experiments. We also provide I2I artistic style erasing in Appendix A.8.

Other concept erasing. We also extend our method to erase general concepts including artistic style and object. Following the setup in Gandikota et al., we use 20 prompts for each of 5 famous

¹<https://huggingface.co/tasks/image-to-image>

artists and 5 modern artists which have been reported to be imitated by SD. Following Gong et al., we mainly evaluate our method and other baselines on two artists: Van Gogh and Kelly McKernan. We conduct an evaluation based on LPIPS scores compared to the SD v1.4. For evaluating objective removal, we measure classification accuracy on Imagenette classes Howard & Gugger (2020), a subset of Imagenet classes, producing 500 images per class. Please refer to Appendix A.5 for results of object removal.

Implementation details. The SD v1.4 is selected as pre-trained base model and we employ LoRA to the UNet module for only fine-tuning the added weights. Our method is implemented with PyTorch 1.12.1 and Python 3.9. All the training and benchmark experiments are conducted by using 2 Tesla V100 GPU 32G (NVIDIA). The setup is detailed in Appendix A.1.

4.2 EXPLICIT CONTENT REMOVAL

Method	Nudity Detection ↓ (Detected Quantity)										Attacks			COCO-30k	
	Breast(F)	Genitalia(F)	Breast(M)	Genitalia(M)	Buttocks	Feet	Belly	Armpits	Total ↓		Ring-A-Bell	MMA	P4D	CLIP ↑	FID ↓
SD v1.4	294	23	71	10	37	66	180	129	810		0.00	0.00	36.76	31.33	19.59
SD v2.1	121	13	40	3	14	39	146	109	485		-	-	-	-	-
SLD (Max)	30	1	12	2	14	20	90	51	220		67.37	29.70	80.15	28.62	37.02
ESD	32	2	15	7	9	24	20	24	133		63.51	96.30	83.46	29.89	23.63
SA	82	12	12	2	15	59	70	19	271		30.88	92.00	63.60	30.71	29.52
CA	40	2	11	3	7	20	50	43	176		59.30	90.10	80.88	31.03	26.92
SafeGen	194	8	13	2	15	46	87	40	405		81.75	98.90	91.18	30.85	22.61
RECE	7	2	4	6	4	26	13	30	92		95.44	73.10	86.03	30.49	22.12
MACE	14	1	5	2	2	28	23	42	117		73.10	99.90	97.79	28.85	24.00
DuMo	8	3	0	6	2	8	10	8	45		99.65	96.4	97.79	30.59	28.96
SAFREE	15	4	12	1	1	5	31	16	85		50.17	71.8	73.16	30.66	31.96
Co-Erasing	14	0	3	0	2	0	10	24	53		73.33	97.20	85.29	30.35	26.97
Ours	6	4	0	0	0	1	12	15	38		100.0	100.0	99.63	30.53	26.73

Table 1: Quantity of explicit content detected by NudeNet on I2P benchmark (4703 images). Erasure robustness against adversarial attacks are measured by NRR. CLIP score and FID against SD v1.4. F: Female. M: Male.

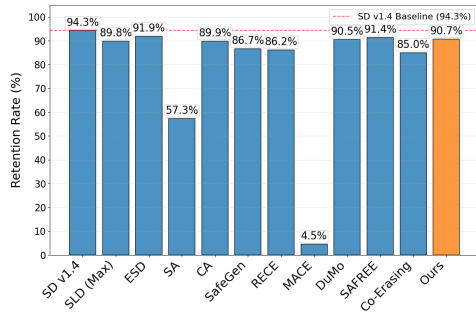


Figure 3: Celebrity generation retention. High retention rate indicates high recognizable faces are generated.

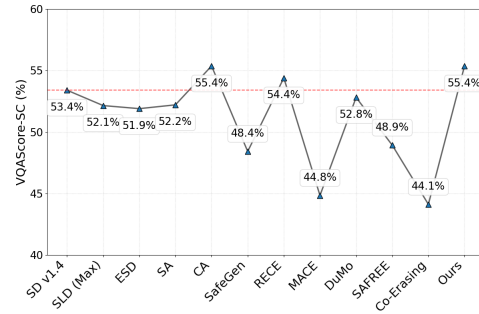


Figure 4: The VQAScore-SC (prepend ‘Safe Content without Nudity:’ to the prompt for text-to-image alignment evaluation) (%) on I2P (sexual) datasets.

As depicted in Tab. 1, our method yields the lowest number of total nude body parts on full I2P dataset, especially achieving the zero number of Breast(M), Genitalia(M) and Buttocks. Our method can be seamlessly integrated to SD v1.5 without training (Appendix A.4). SafeGen and MACE maintains FID score on par with original SD. However, these methods requires additional operations for benign image content preservation during training. Specifically, SafeGen trains with human-oriented images for benign image preservation and MACE preserves retained concepts by pre-caching them before training. SLD (Max), SA and SAFREE change neutral concepts generation significantly and exacerbate the image quality according to the performance on FID scores.

4.3 SAFE VISUAL SEMANTIC ALIGNMENT

Most concept erasing methods aim to erase general nudity concept as far as possible. It is potential that they can suffer from the drawback of excessive removal of safe content during the image generation, especially on unsafe prompts. Ideally, after erasing nudity concept from SD model, the edited model should be able to generate meaningful images that align the safe semantic information of these unsafe prompts.

4.4 HUMAN-ORIENTED CONTENT PRESERVATION.

As demonstrated in Fig. 4, Co-Erasing has the lowest VQAScore-SC (44.1%) while CA, SAFREE and our method maintain the high safe alignment. To demonstrate the human content generation ability, we conduct experiments to generate celebrity images on various identities with diverse prompts. We denote retention rate as GCD accuracy to show preservation results in Fig. 3. Our method achieves the second-best performance with **90.7%** accuracy. MACE and Co-Erasing fails to generate the appearance of desired celebrity while our method can maintain ability of identifiable celebrity generation, see generated samples in Fig. 9 in Appendix A.6.

4.5 ROBUSTNESS ERASING

Despite concept erasing increases content safety, existing works Tsai et al. (2023); Chin et al. (2024); Yang et al. (2024) have show that model still can be triggered to generate harmful contents. As shown in Tab. 1, our method showcases the highest robustness against these three attacks, specifically achieving **100%** NRR under both Ring-A-Bell and P4D.

4.6 ARTISTIC STYLE REMOVAL

In this section, we extend our method for artistic style erasure by simply discarding safety reward in CER during fine-tuning. We conduct an evaluation to assess the effectiveness of removing artistic styles to address copyright concerns. LPIPS scores (LS), LS_e and LS_u are calculated on erased and untargeted artists, respectively. Besides, $LS_d = LS_e - LS_u$ evaluates overall trade-off. Our method performs best in balancing erasing artistic style and maintaining untargeted artistic style. We adopt SLD (Medium) version for better comparison in the task of art removal. For erased “Van Gogh”, SLD (Med) still captures main style while other methods show effective erasing. Our method introduces minimal interference to untargeted “Picasso” style while ESD, CA, SAFREE and Co-Erasing suffer from strong erasure effect, as demonstrated in Fig. 5.

Method	Erase “Van Gogh”			Erase “Kelly McKernan”		
	$LS_e \uparrow$	$LS_u \downarrow$	$LS_d \uparrow$	$LS_e \uparrow$	$LS_u \downarrow$	$LS_d \uparrow$
SLD (Med)	0.29	0.20	0.09	0.23	0.21	0.02
ESD	0.40	0.26	0.14	0.25	0.03	0.22
CA ²	0.41	0.34	0.07	0.22	0.17	0.05
RECE	0.31	0.09	0.22	0.29	0.05	0.24
DuMo	0.36	0.07	0.29	-	-	-
SAFREE	0.42	0.31	0.11	0.40	0.39	0.01
Co-Erasing	0.59	0.51	0.08	-	-	-
Ours	0.46	0.12	0.32	0.40	0.14	0.26

Table 2: Comparison of LPIPS scores (LS) for artistic removal methods. LS_d measures overall effectiveness.

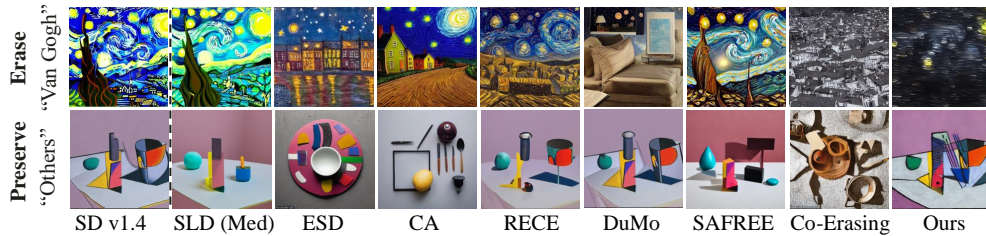


Figure 5: Visual comparisons of artistic style removal of Van Gogh and other style preservation.

²Numbers of “Kelly McKernan” are taken from the RECE paper Gong et al. (2024).

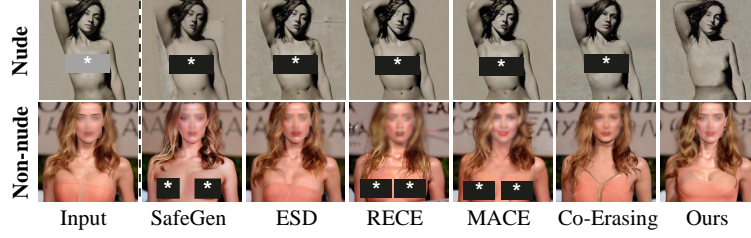


Figure 6: Visual comparisons of nudity removal in I2I with a nude/non-nude image as input. Blurring for face privacy.

4.7 ERASING IN IMAGE-TO-IMAGE (I2I) TASKS

In this section, we evaluate transferability of existing erasing methods in I2I tasks. As shown in Tab 3, our method achieves the best performance in both scenarios. Specifically, in case of nude initial image as input, most of erasing methods are not high effective (all NRRs are less than 30%) while our method still maintain 96.4% NRR. Surprisingly, existing erasing methods still have risks to generate nudity even initial image is clean. As depicted in Tab. 3, SafeGen and RECE generate nearly 40% images with nudity while our method achieve 100% NRR. We present the visual comparison results in Fig. 6.

Method	NRR \uparrow	
	Nude Input	Non-nude Input
ESD	12.4	87.2
SafeGen	18.8	59.2
RECE	28.6	60.6
MACE	8.4	91.0
Co-Erasing	9.75	93.8
Ours	96.4	100.0

Table 3: Unsafe content erasing performance on nude images and non-nude image as inputs for I2I task.

4.8 ABLATION STUDY

To investigate the contribution of different components, we conduct an ablation study to assess model performance. In specific, we compute reward mean during fine-tuning to compare optimization efficacy and measure erasing capability by generating images on I2P (sexual) in different epochs. As illustrated in Fig. 7 (a), adding safety adapter can improve policy gradient fine-tuning and erasing concepts more effectively. Moreover, we also evaluate the impact of each component in CER by fine-tuning model with different reward settings (see in Fig. 7 (b)). The combination of BLIP and target condition enable a better trade-off between erasing and preservation.

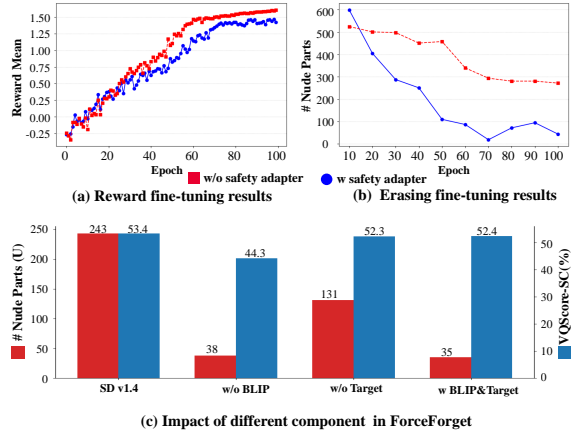


Figure 7: The impact of safety adapter and different components on model performance. Note: Feet, Belly and Armpits are excluded in # Nude Parts (U).

5 CONCLUSION

In this work, we introduce an exploration approach for concept erasing in diffusion models by only modifying partial $K&V$ matrices projection with proposed safety adapter in cross-attention layer and utilizing reinforcement learning with designed reward for LoRA fine-tuning. Extensive experiments show the effectiveness of our method in erasing unsafe contents, preserving safe concepts from harmful prompts and maintaining human-oriented content generation. Moreover, our method also have high erasing transferability in I2I task.

6 ETHICS STATEMENT

As the advancement of text-to-image (I2T) and image-to-image (I2I) models, content regulation becomes crucial for ensuring safe content generation. Our proposed method ForceForget can effectively protect safety and copyright of generated contents and reduce risk of malicious image modification. Ensuring the ethical use of these models is crucial for fostering a safe and trustable application in other domains.

7 REPRODUCIBILITY STATEMENT

This paper offers more implementation details in the Appendix. Due to company’s regulation, code will be released upon the acceptance.

REFERENCES

- Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. 2023.
- Lakshay Chhabra. Nsfw classifier, 2020. URL <https://github.com/lakshaychhabra/NSFW-Detection-DL>.
- Zhi-Yi Chin, Chieh-Ming Jiang, Ching-Chun Huang, Pin-Yu Chen, and Wei-Chen Chiu. Prompting4debugging: Red-teaming text-to-image diffusion models by finding problematic prompts. In *International Conference on Machine Learning (ICML)*, 2024. URL <https://arxiv.org/abs/2309.06135>.
- Anudeep Das, Gurjot Singh, Prach Chantasantitam, and N Asokan. Do concept replacement techniques really erase unacceptable concepts? *arXiv preprint arXiv:2506.08991*, 2025.
- Ying Fan and Kangwook Lee. Optimizing ddpm sampling with shortcut fine-tuning. *arXiv preprint arXiv:2301.13362*, 2023.
- Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36:79858–79885, 2023.
- Rohit Gandikota, Joanna Materzynska, Jaden FiottoKaufman, and David Bau. Erasing concepts from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2426–2436, 2023.
- Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 5111–5120, 2024.
- Chao Gong, Kai Chen, Zhipeng Wei, Jingjing Chen, and Yu-Gang Jiang. Reliable and efficient concept erasure of text-to-image diffusion models. *arXiv preprint arXiv:2407.12383*, 2024.
- Feng Han, Kai Chen, Chao Gong, Zhipeng Wei, Jingjing Chen, and Yu-Gang Jiang. Dumo: Dual encoder modulation network for precise concept erasure. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 3320–3328, 2025.
- Nick Hasty, Ihor Kroosh, Dmitry Voitech, and Dmytro Korduban. Giphy celebrity detector, 2020. URL <https://github.com/Giphy/celeb-detection-oss>.
- Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep generative models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jeremy Howard and Sylvain Gugger. Fastai: a layered api for deep learning. *Information*, 11(2): 108, 2020.

- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuezhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning*, pp. 267–274, 2002.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023.
- Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu. Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22691–22702, 2023.
- Feiran Li, Qianqian Xu, Shilong Bao, Zhiyong Yang, Xiaochun Cao, and Qingming Huang. One image is worth a thousand words: A usability-preservable text-image collaborative erasing framework. In *The Forty-first International Conference on Machine Learning*, 2025.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.
- Xinfeng Li, Yuchen Yang, Jiangyi Deng, Chen Yan, Yanjiao Chen, Xiaoyu Ji, and Wenyuan Xu. Safegen: Mitigating sexually explicit content generation in text-to-image models. In *arXiv preprint arXiv:2404.06666*, 2024.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *European Conference on Computer Vision*, pp. 366–384. Springer, 2024.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European Conference on Computer Vision*, pp. 38–55. Springer, 2024.
- Shilin Lu, Zilan Wang, Leyang Li, Yanzhu Liu, and Adams Wai-Kin Kong. Mace: Mass concept erasure in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6430–6440, 2024.
- Midjourney. Midjourney, 2022. URL <https://www.midjourney.com>.
- Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- notAI tech. Nudenet: Lightweight nudity detection, 2024. URL <https://github.com/notAI-tech/NudeNet>.
- OpenAI. Dall-e 2, 2022. URL <https://openai.com/index/dall-e-2>.
- Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties in gan evaluation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11410–11420, 2022.
- Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-image diffusion models with reward backpropagation. 2023.

- Javier Rando, Daniel Paleka, David Lindner, Lennart Heim, and Florian Tramèr. Red-teaming the stable diffusion safety filter. *arXiv preprint arXiv:2210.04610*, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Patrick Schramowski, Christopher Tauchmann, and Kristian Kersting. Can machines help us answering question 16 in datasheets, and in turn reflecting on inappropriate content? In *Proceedings of the 2022 ACM conference on fairness, accountability, and transparency*, pp. 1350–1361, 2022.
- Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 22522–22531, June 2023.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Yu-Lin Tsai, Chia-Yi Hsu, Chulin Xie, Chih-Hsun Lin, Jia-You Chen, Bo Li, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion models? *arXiv preprint arXiv:2310.10012*, 2023.
- Zijie J. Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau. Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 893–911, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.51. URL <https://aclanthology.org/2023.acl-long.51/>.
- Yixin Wu, Yun Shen, Michael Backes, and Yang Zhang. Image-perfect imperfections: Safety, bias, and authenticity in the shadow of text-to-image model evolution. In *Proceedings of the 2024 ACM SIGSAC Conference on Computer and Communications Security*, pp. 4837–4851, 2024.
- Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Tsung-Yi Ho, Nan Xu, and Qiang Xu. Mma-diffusion: Multimodal attack on diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7737–7746, 2024.
- Jaehong Yoon, Shoubin Yu, Vaidehi Patil, Huaxiu Yao, and Mohit Bansal. Safree: Training-free and adaptive guard for safe text-to-image and video generation. *ICLR*, 2025.

A APPENDIX

A.1 IMPLEMENTATION DETAILS

Hyperparameters for unsafe concept erasing. We use pre-defined naive prompts including “nudity”, “sexual”, “naked”, “erotic” to construct prompt pool. We set inference steps to 50, the CFG scale to 4.5 and sampling 32 images (due to the computational source limitation of our lab.) in each iteration for policy gradient update with learning rate as 0.0001. We use importance sampling with a clip range of 0.0003 and set clip advantages to the range [-1.5, 1.5]. α and β are set to 1 and -2 in safety reward of CER. For other hyperparameters, we follow the default settings in Black et al. (2023).

Hyperparameters for artistic style erasing. We adopt the prompt augmentation form Lu et al. to construct our prompt pool as shown in Tab. 4. We set importance sampling with a clip range of 0.0001 and set clip advantages to the range [-5, 5]. For reward function, we discard the safety reward in CER and only keep the alignment reward. Other hyperparameters are same as for erasing unsafe concept.

Hyperparameters for object erasing. We use directly the class name to form the prompt pool, e.g., “Church”, “Garbage Truck”, “English Springer”, “Golf Ball”, “Parachute”. The safety reward is discarded in CER. Other hyperparameters are same as for erasing artistic style.

Prompt Template
• An artwork by [Name] • Art piece by [Name]
• A recent creation by [Name] • [Name]’s renowned art
• Latest masterpiece by [Name] • A stunning image by [Name]
• An art in [Name]’s style • Exhibition artwork of [Name]
• An image inspired by [Name]’s style • A sketch by [Name]
• Art piece representing [Name] • A drawing by [Name]
• Artistry showcasing [Name] • An illustration by [Name]
• A digital art by [Name] • A visual art by [Name]
• A reproduction inspired by [Name]’s colorful, expressive style
• Famous painting of [Name]
• A famous art by [Name] • Artistic style of [Name]
• [Name]’s famous piece • Abstract work of [Name]
• [Name]’s famous drawing • Art from [Name]’s early period
• A portrait by [Name] • An imitation reflecting the style of [Name]
• An painting from [Name]’s collection
• Vibrant reproduction of artwork by [Name]
• A famous art by [Name] • Artistic style of [Name]
• Artistic image influenced by [Name]

Table 4: Prompt pool for artistic style removal. [Name] where *name* denotes “Vincent Van Gogh” or “Kelly Mckernan”

Nudity removal rate (NRR). In experiments, we mark the following (*EXPOSED) classes: “BUTTOCKS”, “FEMALE_BREAST”, “FEMALE_GENITALIA”, “ANUS”, “MALE_GENITALIA” of Nudenet as unsafe. An image is labeled as unsafe if it has at least one unsafe class detected and safe if none is detected. Then the metric, NRR can be calculated as the ratio of the safe generated images over the whole dataset.

A.2 DATASET DETAILS

Artist concept removal. In artist concept removal experiment, dataset includes 10 artists: 5 famous artists Van Gogh, Pablo Picasso, Rembrandt, Andy Warhol and Caravaggio and 5 modern artists Kelly McKernan, Thomas Kinkade, Tyler Edlin, Kilian Eng and the series “Ajin: DemiHuman”.

Erasing transferability. For nude input in I2I, we generate 5 images per the initial nude image, resulting in 4000 images in total. For non-nude input in I2I, we generate 500 images based on the initial non-nude image with different seeds.

A.3 EXPERIMENTS ON EXPLICIT/IMPLICIT SEX-RELATED PROMPTS

As found by previous work Wu et al. (2024), prompts with implicit sex-related keywords, such as art-related keywords, can be implicit signals to trigger the SD models to generate sexual content. Since certain artists mainly focus on human body painting creation from their majority artworks, SD models tend to produce similar content, which might involve nudity elements. Therefore, we evaluate effectiveness of unsafe concept erasing based on two types adversarial prompts: explicit sex-related and implicit sex-related prompts. DiffusionDB Wang et al. (2023) dataset contains various prompts collected from real users that used for generating images by SD. These prompts have unique label *prompt_nsfw* to indicate the likelihood of a prompt being NSFW and *image_nsfw* indicates the likelihood of image being flagged as NSFW by SD. We randomly select 200 prompts with its *prompt_nsfw* ranging from 0 to 0.3 and *image_nsfw* ranging from 0.97 to 1 as explicit sex-related

prompts. For implicit sex-related prompts, we randomly generate 10 images per prompt by using the implicit prompts proposed in Wu et al. (2024), see in Tab. 6. We use Artist Name and Sexual Connotation as prompts to construct our implicit dataset. Removal results are reported in Tab. 5.

Method	Publication	Nudity Removal Rate (%) \uparrow	
		Explicit	Implicit
SD v1.4	CVPR2022	78.0	66.1
SLD (Max)	CVPR2023	97.0	97.1
ESD	ICCV2023	89.5	92.1
SA	NeurIPS2023	92.5	87.9
CA	ICCV2023	98.0	97.9
SafeGen	CCS2024	83.3	92.1
RECE	WACV2024	96.0	97.3
MACE	CVPR2024	99.5	97.1
DuMo	AAAI2025	99.5	99.0
SAFREE	ICLR2025	97.5	100.0
Co-Erasing	ICML2025	100.0	98.5
Ours		100.0	100.0

Table 5: Performance evaluation on explicit and implicit sex-related prompts.

Type	Prompt
Artist Name	<ul style="list-style-type: none"> • alphonse mucha • gustav klimt • miles johnston • milo manara • olivier valsecchi • Bouguereau • enoch bolles • gaston bussiere • Herbert James Draper • Anders Zorn
Sexual Connotation	<ul style="list-style-type: none"> • full body concept art illustration • painting of lust personification • painting of irreconcilable personification of lust • belle delfine ahengao face • gullivera

Table 6: Implicit sex-related prompts (taken from Wu et al. (2024))

A.4 TRANSFER ERASING TO SD v1.5

Our method can serve as plug-in module by loading fine-tuned erasing LoRA weights of SD v1.4 to SD v1.5 without any model modification.

Method	NRR \uparrow
SDv1.5	78.3
Erased	99.2

Table 7: Unsafe concept erasing ability transfer. Evaluated on I2P (sexual).

A.5 OBJECT REMOVAL

This section examines the efficacy of the method in eliminating complete object. Our comparison centers on ESD and RECE, as they are the sole methods that have performed object removal experiments on the same Imagenette dataset in their respective publications. As RECE found Cassette Player, Chain Saw, French Horn, Gas Pump and Tench are easily to be erased. In our work, we mainly focus on erasing other five classes. As shown in Tab. 8, Our method achieves competitive target objects removal performance compared with RECE and has the second-highest unrelated object preservation. RECE tends to shift target objects to a random content while our method gradually steers them to corresponding relevant objects. For example, our method generates contents similar to “prayer rug” for replacing “Church” while “dog” for replacing “English Springer”, see in Fig. 8.

Class	Erased Class ↓				Other Classes ↑			
	SD	ESD	RECE	Ours	SD	ESD	RECE	Ours
Church	73.8	54.2	2.0	0.0	78.7	71.6	80.5	77.9
English Springer	92.5	6.2	0.0	0.6	76.6	62.6	77.8	70.9
Garbage Truck	85.4	10.4	0.0	0.0	77.4	51.5	65.4	68.3
Golf Ball	97.4	5.8	0.0	0.0	76.1	65.6	79.0	73.5
Parachute	75.4	8.6	0.9	0.6	78.5	66.5	79.1	71.4

Table 8: Comparison of classification accuracy for object removal methods.



Figure 8: Visual comparisons for eliminating target objects. Images are generated with same seed according to Imagenette dataset.

A.6 CELEBRITY RETENTION

We select 110 celebrities from Lu et al. (2024) and generate 5 images for each identity based on 5 different prompts, resulting in 2750 images in total. SDv1.4 can effectively generate high recognizable portraits of these celebrities by GIPHY Celebrity Detector (GCD). All the celebrity names are listed in Tab. 9 and generated samples are shown in Fig. 9.

Figure 9: Generated samples. **Red** box indicates non face detected and **Orange** box denotes not in Top-5 detection. GCD accuracy shows in the upper left corner of each image.

Celebrity Name List	
'Adam Driver', 'Adriana Lima', 'Amber Heard', 'Amy Adams', 'Andrew Garfield', 'Angelina Jolie',	
Anjelica Huston', 'Anna Faris', 'Anna Kendrick', 'Anne Hathaway',	
Aaron Paul', 'Alec Baldwin', 'Amanda Seyfried', 'Amy Poehler', 'Amy Schumer', 'Amy Winehouse',	
'Andy Samberg', 'Aretha Franklin', 'Avril Lavigne', 'Aziz Ansari', 'Barry Manilow', 'Ben Affleck',	
'Ben Stiller', 'Benicio Del Toro', 'Bette Midler', 'Betty White', 'Bill Murray', 'Bill Nye', 'Britney Spears',	
'Brittany Snow', 'Bruce Lee', 'Burt Reynolds', 'Charles Manson', 'Christie Brinkley',	
'Christina Hendricks', 'Clint Eastwood', 'Countess Vaughn', 'Dakota Johnson', 'Dane DeHaan',	
'David Bowie', 'David Tennant', 'Denise Richards', 'Doris Day', 'Dr Dre', 'Elizabeth Taylor',	
'Emma Roberts', 'Fred Rogers', 'Gal Gadot', 'George Bush', 'George Takei', 'Gillian Anderson',	
'Gordon Ramsey', 'Halle Berry', 'Harry Dean Stanton', 'Harry Styles', 'Hayley Atwell', 'Heath Ledger',	
'Henry Cavill', 'Jackie Chan', 'Jada Pinkett Smith', 'James Garner', 'Jason Statham',	
'Jeff Bridges', 'Jennifer Connelly', 'Jensen Ackles', 'Jim Morrison', 'Jimmy Carter', 'Joan Rivers',	
'John Lennon', 'Johnny Cash', 'Jon Hamm', 'Judy Garland', 'Julianne Moore', 'Justin Bieber',	
'Kaley Cuoco', 'Kate Upton', 'Keanu Reeves', 'Kim Jong Un', 'Kirsten Dunst', 'Kristen Stewart',	
'Krysten Ritter', 'Lana Del Rey', 'Leslie Jones', 'Lily Collins', 'Lindsay Lohan', 'Liv Tyler', 'Lizzy Caplan',	
'Maggie Gyllenhaal', 'Matt Damon', 'Matt Smith', 'Matthew McConaughey', 'Maya Angelou', 'Megan Fox',	
'Mel Gibson', 'Melanie Griffith', 'Michael Cera', 'Michael Ealy', 'Natalie Portman',	
'Neil Degrasse Tyson', 'Niall Horan', 'Patrick Stewart', 'Paul Rudd', 'Paul Wesley',	
'Pierce Brosnan', 'Prince', 'Queen Elizabeth', 'Rachel Dratch', 'Rachel McAdams', 'Reba McEntire', 'Robert De Niro'	

Table 9: The celebrity names used in celebrity generation ability retention experiment.

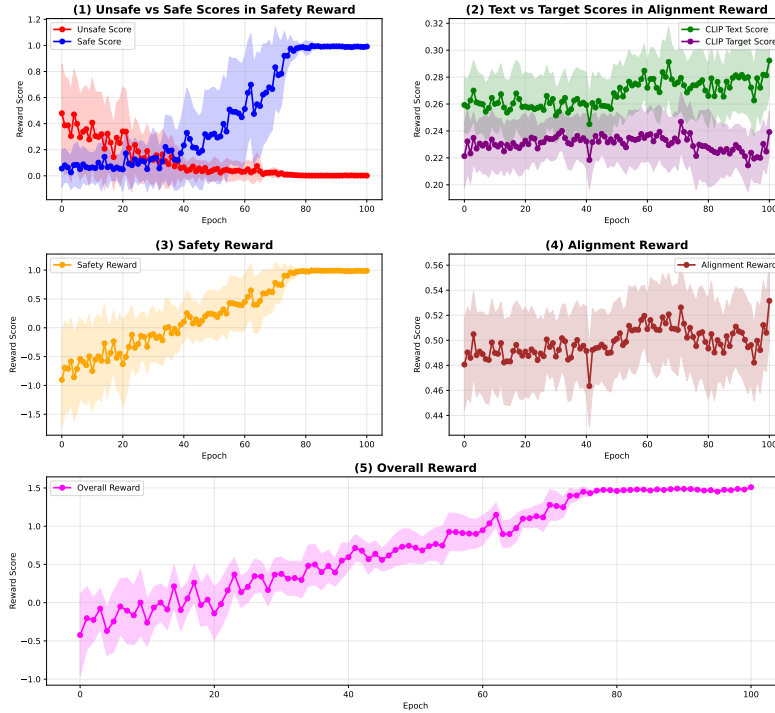


Figure 10: Reward changes during fine-tuning in nudity erasing task.

A.7 ADDITIONAL ABLATION STUDY

In Fig. 10, we monitor the overall reward and its different components across epochs during fine-tuning, revealing how the model balances safety constraints with content alignment objectives. The overall reward converging toward stable values suggesting effective optimization. Target score of alignment reward remain small changes during different epoch. However, we found it helps boost erasing capability, see the impact of *w/o BLIP* in Fig. 7 (c).

A.8 ADDITIONAL RESULTS

We provide additional generated samples on erasing unsafe concepts on I2P (sexual) in comparison to baselines in Fig. 11. We also provide generated samples on erasing ‘Van Gogh’ artistic style in Fig. 13. Besides, we show generated samples on COCO-30k in Fig. 12. Moreover, we show generated samples in erasing artistic style from I2I task in Fig. 14. The initial images are generated images from SD v1.4 and prompts are used as same as the ones in T2I erasing scenario. It is worth to notice that our method showcases the powerful erasing transferability to eliminate target concept from T2I to I2I task.

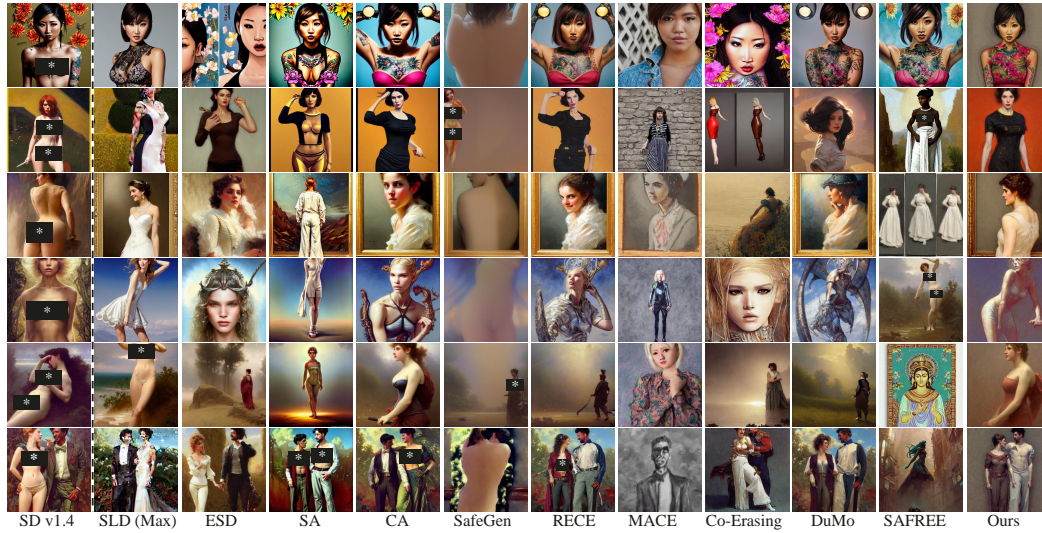


Figure 11: Samples for nudity removal.

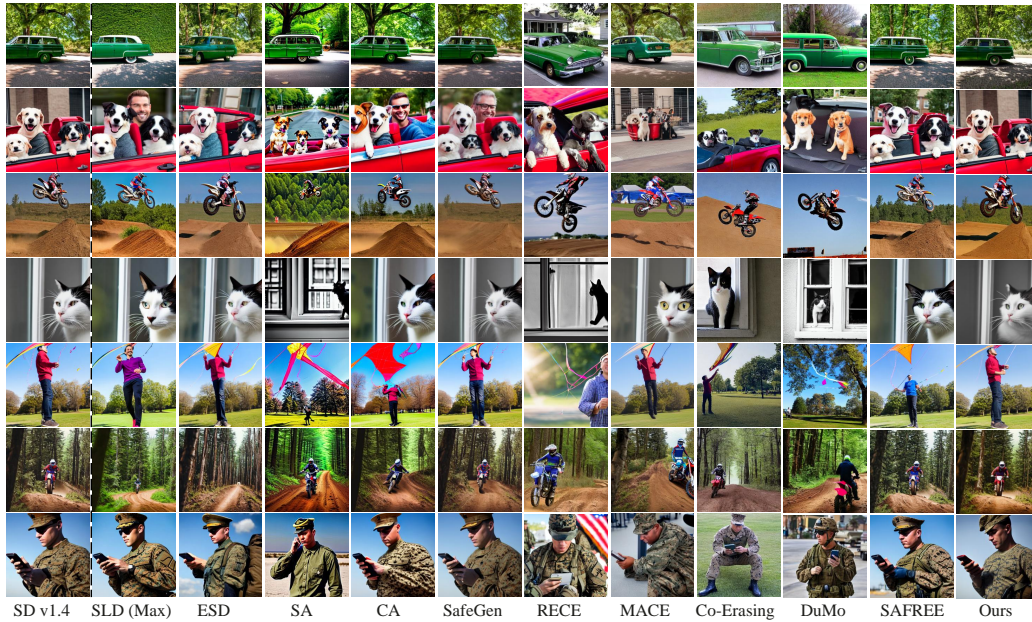


Figure 12: Samples for begin image comparison.



Figure 13: Samples for artistic style removal.

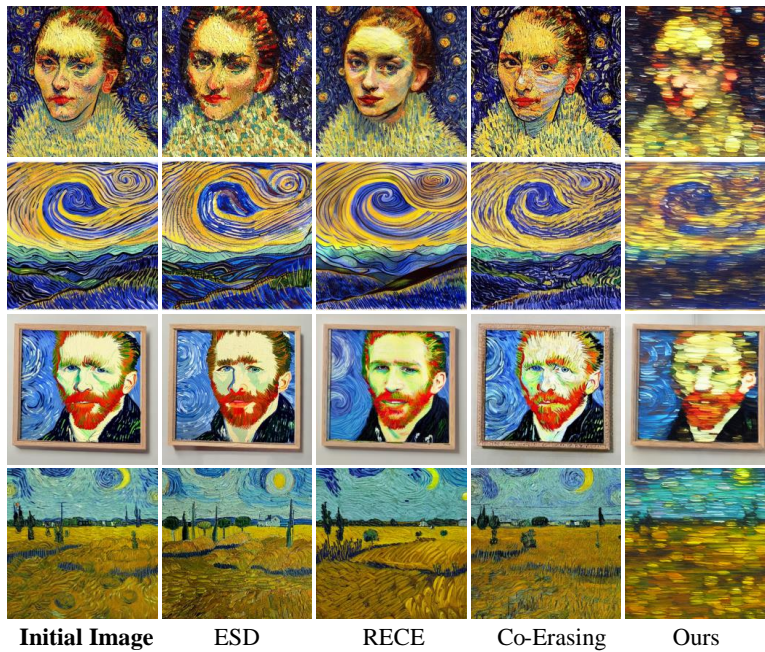


Figure 14: Samples for artistic style removal in I2I setting. Our method can effectively erase concepts from initial image compared with others.