

Set You Straight: Auto-Steering Denoising Trajectories to Sidestep Unwanted Concepts

Leyang Li^{1,*} Shilin Lu^{1,*} Yan Ren¹ Adams Wai-Kin Kong¹

¹Nanyang Technological University, Singapore

{lile0005, shilin002}@e.ntu.edu.sg, nomatterhowlong@gmail.com, adamskong@ntu.edu.sg

Abstract

*Ensuring the ethical deployment of text-to-image models requires effective techniques to prevent the generation of harmful or inappropriate content. While concept erasure methods offer a promising solution, existing finetuning-based approaches suffer from notable limitations. Anchor-free methods risk disrupting sampling trajectories, leading to visual artifacts, while anchor-based methods rely on the heuristic selection of anchor concepts. To overcome these shortcomings, we introduce a finetuning framework, dubbed **ANT**, which **A**utomatically guides **d**eNoising **T**rajectories to avoid unwanted concepts. **ANT** is built on a key insight: reversing the condition direction of classifier-free guidance during mid-to-late denoising stages enables precise content modification without sacrificing early-stage structural integrity. This inspires a trajectory-aware objective that preserves the integrity of the early-stage score function field—which steers samples toward the natural image manifold—without relying on heuristic anchor concept selection. For single-concept erasure, we propose an augmentation-enhanced weight saliency map to precisely identify the critical parameters that most significantly contribute to the unwanted concept, enabling more thorough and efficient erasure. For multi-concept erasure, our objective function offers a versatile plug-and-play solution that significantly boosts performance. Extensive experiments demonstrate that **ANT** achieves state-of-the-art results in both single and multi-concept erasure, delivering high-quality, safe outputs without compromising the generative fidelity. Code is available at <https://github.com/lileyang1210/ANT>.*

1. Introduction

Concept erasure in text-to-image (T2I) models [5, 11, 60, 67, 69, 71, 86] addresses the critical challenge of preventing the generation of harmful or inappropriate visual content, such as violent, explicit, copyright-infringing, or offensive imagery. Current methods for concept erasure can

be broadly categorized into two types: **(1) finetuning-based methods** [19, 55, 57], which directly modify model parameters, and **(2) finetuning-free methods** [33, 58, 73], which aim to influence model outputs without parameter updates. However, finetuning-free methods are vulnerable to bypassing when the source code is openly available, thus making finetuning-based methods more effective and secure for publicly accessible models.

Finetuning-based methods remove undesirable data modes by altering the predicted score function field—essentially, modifying the gradient directions that samples follow during the denoising process—to avoid converging toward undesirable image distributions. As a result of finetuning, the predicted score function no longer accurately reflects the true gradient direction in data space that would further increase likelihood. The main difference among finetuning-based techniques lies in how the conditional score function is modified, which can be broadly divided into **anchor-free** and **anchor-based** approaches.

Anchor-free methods [3, 7, 8, 10, 19, 21, 24, 28, 30, 32, 35, 37–39, 53, 57–59, 72, 73, 80, 81, 84, 85, 87, 88, 95, 100] often design a loss to adjust the conditional score function throughout the denoising process, encouraging samples to move away from unwanted image manifolds without explicitly specifying a target manifold (see Figure 1(b)). However, this approach can disrupt the sampling trajectories toward natural image manifolds. As shown in Figure 1(a), diffusion models typically first guide samples from Gaussian noise toward the manifold of natural images to establish a plausible layout, and then progressively refine the details during the mid-to-late denoising steps [41, 55]. By solely emphasizing the movement away from unwanted manifolds, anchor-free methods risk causing samples to deviate from the natural image manifold early on, potentially resulting in generated images with visual artifacts or unintended content (see the second row of Figure 2).

Anchor-based methods [2–4, 6, 9, 15, 16, 20, 23, 26, 31, 36, 40, 43, 48, 49, 52, 55, 62, 74, 78, 82, 89, 90, 99, 102], on the other hand, typically utilize a loss designed to leverage benign anchor concepts by aligning the predicted condi-

* Equal contribution

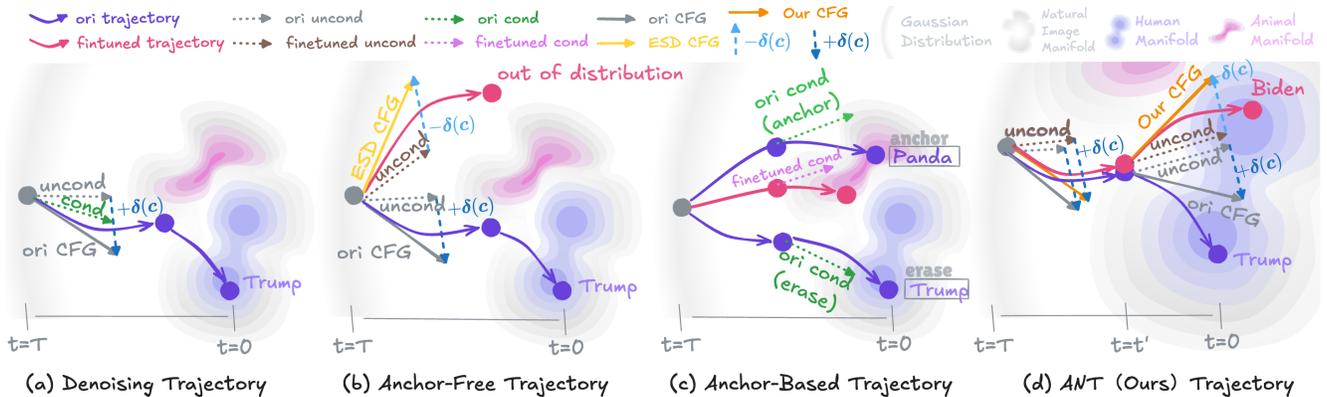


Figure 1. Geometric perspective on concept erasure in diffusion models. **(a) Conventional Denoising Trajectory.** A high-dimensional Gaussian sample, starting on a large sphere, converges to the human data manifold via classifier-free guidance (CFG). **(b) Anchor-Free Finetuned Trajectory.** Finetuning often modifies the orientation of the predicted conditional score functions so that they direct away from the unwanted concept manifold. This results in a condition direction $\delta(c) = \epsilon_{\theta}(z_t, t, c) - \epsilon_{\theta}(z_t, t)$ nearly opposite to that of the original model, making the trajectory more likely to produce out-of-distribution samples. Note that, in the absence of an unconditional constraint, modifications to the conditional output also affect the unconditional output due to shared model parameters. **(c) Anchor-Based Finetuned Trajectory.** The model is finetuned so that the predicted score functions (or keys & values) for the unwanted concept align with those of the original model conditioned on a benign anchor, ensuring final samples lie on the anchor manifold, though not necessarily at the highest-probability mode. **(d) Our Trajectory (ANT).** In the early stage (when $t > t'$), the conditional score functions remain directed toward the natural data mode, keeping the finetuned model aligned with the original. When $t < t'$, they are finetuned to point away from the unwanted concept manifold. ANT encourages that unconditional score functions remain unchanged throughout all stages.

tional score functions (or keys & values) for unwanted concepts with those associated with anchor concepts (see Figure 1(c)). By aligning score functions of unwanted concepts with those of anchor concepts, these methods ensure that samples conditioned on unwanted concepts ultimately converge towards images depicting the anchor concepts. Thus, anchor-based approaches are not merely repelling samples from undesired modes. Nevertheless, the effectiveness of these methods critically depends on the proper selection of anchor concepts. As demonstrated in the third row of Figure 2, some seemingly reasonable anchor concept choices can reduce the quality of images generated when conditioned on erased concepts. Currently, selecting effective anchor concepts remains largely heuristic, lacking systematic guidelines.

Motivated by these limitations, we propose a trajectory-aware finetuning framework, termed ANT, which Automatically guides deNoising Trajectories to avoid unwanted concepts. This approach achieves its goal without negatively affecting early-stage score function fields or relying on heuristic anchor concept selection. Specifically, we discovered that reversing the condition direction of classifier-free guidance (CFG) [29] during the mid-to-late denoising stage enables modification of detailed content while preserving the fundamental structure of the generated image. This finding inspires us to develop a trajectory-aware objective function that preserves the early-stage score function, steering samples toward the natural image

manifold, while eliminating the need for anchor concepts (Figure 1(d)). This approach enables more effective erasure of undesired concepts while better preserving those that are unrelated. In the context of single-concept erasure, we introduce an augmentation-enhanced weight saliency map that accurately identifies the key parameters most responsible for generating a specific concept. Moreover, our loss function is fully compatible with existing multi-concept erasure frameworks, offering a flexible plug-and-play solution, and elevates the performance to a new state-of-the-art (SOTA) level. Our experimental results demonstrate that our method achieves SOTA performance in both single and multi-concept erasure settings. Our contributions are summarized as follows:

1. We offer a geometric perspective on concept erasure and an insight that reversing the condition direction of classifier-free guidance during the mid-to-late denoising stages enables precise content modification while preserving early-stage structural integrity, thus benefiting the erasure community in advancing algorithm designs.
2. We propose a trajectory-aware finetuning framework, which encourages the model to reorient its denoising trajectories during the mid-to-late stages while keeping the early-stage trajectories largely unchanged. This approach enables more thorough erasure of unwanted concepts and better preservation of unrelated ones.
3. We introduce an augmentation-enhanced weight saliency

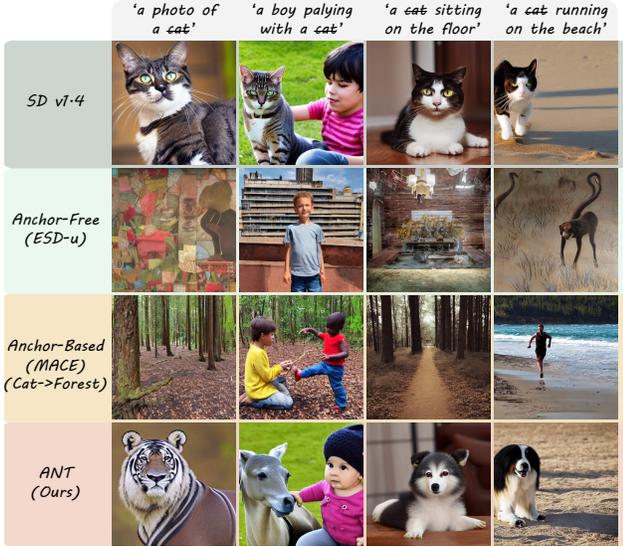


Figure 2. Generation results of different concept erasure methods conditioned on the concept “cat”. The anchor-free method (ESD) often produces images with visual artifacts or content that is out of distribution. The anchor-based method (MACE), which maps “cat” to “forest”, performs reasonably well in simple contexts but results in unnatural or incoherent outputs in more complex scenarios. In contrast, our trajectory-aware method (ANT) effectively removes the target concept while preserving the overall structure and contextual integrity of the generated images.

map that precisely identifies the key parameters most responsible for generating the undesired concept, thereby enabling more effective and efficient erasure.

4. The proposed objective function substantially enhances the performance of existing multi-concept erasure frameworks, achieving SOTA results in both single- and multi-concept erasure settings.

2. Related Work

In this section, we review prior work on concept erasure in diffusion models, with a particular focus on the critical trade-off between erasure and preservation, which is most pertinent to our study. Additional discussions on other dimensions of concept erasure (e.g., finetuning efficiency, scalability, and robustness to adversarial prompts) are provided in Appendix.

The investigation of concept erasure within diffusion models has been pioneered by several foundational studies, establishing the groundwork for this burgeoning domain. SLD [73] introduces an inference-time guidance technique to suppress undesired concepts without modifying the model’s parameters, offering a non-invasive yet effective approach. In contrast, ESD [19] employs direct parameter editing through negative guidance, achieving perma-

nent concept removal. FMN [99] builds upon this trajectory by proposing a lightweight method that manipulates attention mechanisms to enhance computational efficiency. Meanwhile, AC [40] presents a finetuning framework that aligns the score function of an unwanted concept with that of an anchor concept, delivering an alternative strategy for concept ablation.

As concept erasure techniques have matured, the research community has increasingly emphasized the dual objectives of effectively eliminating target concepts while preserving the integrity of unrelated concepts during the finetuning process. Numerous studies [3, 4, 6–8, 16, 17, 21, 22, 24, 26, 28, 32, 37, 39, 48, 52, 55, 57, 58, 72, 74, 80–82, 84, 85, 87, 88, 90, 95, 102] highlight the necessity of maintaining balanced model performance across both targeted and non-targeted concepts. However, a critical limitation of these approaches lies in their insufficient attention to the impacts of finetuning on the early-stage score function. This oversight can lead to a divergence between the predicted score function and the true score function, i.e., the gradient direction in data space that maximizes likelihood. As a result, the generated samples may fail to converge toward the natural image manifold, ultimately degrading the quality and reliability of the outputs. Our work seeks to bridge this gap by explicitly addressing the preservation of the early-stage score function, ensuring both effective concept erasure and high-fidelity generation.

3. Method

We propose ANT, a framework designed to erase specific concepts from pretrained text-to-image diffusion models. Our approach addresses key challenges by eliminating the negative impacts on early-stage score function fields and removing the dependency on heuristic methods for anchor concept selection. The framework requires only two inputs: a pretrained diffusion model and a set of target phrases representing the concepts to be erased. The output is a finetuned model that no longer generates images depicting the unwanted concepts.

3.1. Insights into the Denoising Process

We thoroughly investigated the denoising process in diffusion models and found that applying CFG during the early sampling stage (when $t' < t < T$), and then reversing the CFG’s condition direction term during the mid-to-late sampling stage (when $0 < t < t'$, as shown in Eq. (1)), allows for altering detailed content while preserving the fundamental structure of the image. In other words, the sample avoids converging toward specific unwanted concepts yet remains

within the natural image manifold.

$$\epsilon_{\theta}^{\text{cfg}}(\mathbf{z}_t, t, \mathbf{c}) = \epsilon_{\theta}(\mathbf{z}_t, t) + s \cdot \text{sgn}(t - t') \cdot \delta(\mathbf{c}), \quad (1)$$

$$\text{sgn}(t - t') = \begin{cases} -1, & \text{if } t \leq t' \\ 1, & \text{if } t > t' \end{cases} \quad (2)$$

where the terms $\epsilon_{\theta}^{\text{cfg}}(\mathbf{z}_t, t, \mathbf{c})$, $\epsilon_{\theta}(\mathbf{z}_t, t, \mathbf{c})$, and $\epsilon_{\theta}(\mathbf{z}_t, t)$ denote the classifier-free guidance output, the conditional prediction, and the unconditional prediction, respectively. The difference $\delta(\mathbf{c}) = \epsilon_{\theta}(\mathbf{z}_t, t, \mathbf{c}) - \epsilon_{\theta}(\mathbf{z}_t, t)$ defines the condition direction. t' is a key parameter used to determine the timestep at which the condition direction should be reversed.

As shown in Figure 3(c), if t' is appropriately selected, this approach allows for the targeted removal of specific attributes or details (e.g., occupation, gender, or age) while preserving the naturalness of the generated images. This is because, during the early stage of denoising, the samples follow the correct score function and are guided onto a plausible data manifold. In the later stages, the guidance steers the samples away from certain modes within that manifold. For instance, in Figure 3(c), the occupation changes from doctor to model, gender shifts from male to female, and age transitions from both old and young to middle-aged, all while staying within the human data manifold.

However, if t' is set too early, the early-stage score function will be significantly altered, leading to a loss of the image’s structural integrity (see Figure 3(d)). On the other hand, if t' is set too late, the samples will have already entered the concept-specific mode, and modifications to the late-stage score function will only affect fine details (see Figure 3(b)).

3.2. Trajectory-Aware Loss Function

Inspired by this finding, we aim to preserve the integrity of the early-stage score function field—which guides samples toward the appropriate natural manifold—by introducing constraints during finetuning. Adjustments will be limited exclusively to the mid-to-late stage score function field. This approach ensures that even when the finetuned model is conditioned on the removed concept, the samples can still converge to the appropriate manifold. Specifically, we propose the following finetuning objective:

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{\text{preserve}} + \lambda_1 \cdot \mathcal{L}_{\text{erase}} + \lambda_2 \cdot \mathcal{L}_{\text{uncond-early}} + \lambda_3 \cdot \mathcal{L}_{\text{uncond-late}} \\ &= \mathbb{E}_{\mathbf{z}_{t_1}, \mathbf{c}, t_1 \sim U(t', T)} \left[\|\epsilon_{\theta}(\mathbf{z}_{t_1}, t_1, \mathbf{c}) - \text{sg}[\epsilon_{\theta^*}(\mathbf{z}_{t_1}, t_1) + \eta \delta(\mathbf{c})]\|_2^2 \right] \\ &+ \lambda_1 \mathbb{E}_{\mathbf{z}_{t_2}, \mathbf{c}, t_2 \sim U(0, t')} \left[\|\epsilon_{\theta}(\mathbf{z}_{t_2}, t_2, \mathbf{c}) - \text{sg}[\epsilon_{\theta^*}(\mathbf{z}_{t_2}, t_2) - \eta \delta(\mathbf{c})]\|_2^2 \right] \\ &+ \lambda_2 \mathbb{E}_{\mathbf{z}_{t_1}, \mathbf{c}, t_1 \sim U(t', T)} \left[\|\epsilon_{\theta}(\mathbf{z}_{t_1}, t_1) - \text{sg}[\epsilon_{\theta^*}(\mathbf{z}_{t_1}, t_1)]\|_2^2 \right] \\ &+ \lambda_3 \mathbb{E}_{\mathbf{z}_{t_2}, \mathbf{c}, t_2 \sim U(0, t')} \left[\|\epsilon_{\theta}(\mathbf{z}_{t_2}, t_2) - \text{sg}[\epsilon_{\theta^*}(\mathbf{z}_{t_2}, t_2)]\|_2^2 \right], \end{aligned} \quad (3)$$

where θ represents the parameters undergoing finetuning, while θ^* denotes the original, frozen parameters. The no-



Figure 3. Effect of condition direction reversal at different timesteps. Each column represents a distinct semantic condition, and each row shows generated outputs under varying reversal strategies. (a) displays originally generated images using a diffusion process (timestep 50→1). (b)–(d) show results when the condition direction $\delta(\mathbf{c}) = \epsilon_{\theta}(\mathbf{z}_t, t, \mathbf{c}) - \epsilon_{\theta}(\mathbf{z}_t, t)$ is reversed at different timesteps (25, 35, and 45). With a proper t' , specific attributes can be removed while preserving image naturalness. If t' is too early, structural integrity is lost; if too late, only fine details are affected.

tation $\text{sg}[\cdot]$ indicates the stop-gradient operation. Timesteps t_1 and t_2 are sampled independently from uniform distributions $U(t', T)$ and $U(0, t')$, respectively, with t' being a predefined hyperparameter. Additionally, \mathbf{z}_{t_1} and \mathbf{z}_{t_2} represent the corresponding noisy latent image variables at these timesteps, and η denotes a hyperparameter. *Notably, two timesteps are sampled during each gradient update iteration to effectively balance the gradients associated with concept erasure and the preservation of unrelated concepts.*

Early-stage preservation. The first term $\mathcal{L}_{\text{preserve}}$ ensures that, during the early stage (when $t > t'$), the predicted conditional score function consistently points toward the natural data mode. This preserves the integrity of the early stage score function field. Consequently, when sampling with the finetuned model conditioned on the erased concept, the generated samples can smoothly transition to the natural image manifold.

Mid-to-late-stage erasure. The second term $\mathcal{L}_{\text{erase}}$ emphasizes that at later stage (when $t < t'$), the predicted conditional score function should actively guide samples away from undesirable modes. It differs from the ESD loss [19] in

that the second term is applied exclusively at later timesteps ($t < t'$), whereas the ESD loss spans all timesteps. Including early timesteps in the ESD loss can unintentionally alter the early-stage score function field, frequently causing samples to be incorrectly guided and thereby failing to converge onto the appropriate manifold. To further explore this issue, we conducted an experiment restricting the application of this second loss term solely to mid-to-late denoising steps, specifically aiming to avoid negatively impacting the early-stage score function field. However, even under this restricted condition, the early-stage score function field was still adversely affected, resulting in suboptimal performance (see the ablation study in Table 2). We hypothesize that this outcome arises primarily due to the shared model parameters across all timesteps within the diffusion process.

Unconditional score function preservation. Since the unconditional score function $\epsilon_\theta(z_t, t)$ represents the general direction toward the approximate center of all data modes, modifying it can influence multiple concepts, as demonstrated by our ablation study. Specifically, Table 2 shows that removing 100 celebrity concepts without incorporating unconditional loss terms negatively impacts the preservation of other celebrity concepts. To address this issue, we introduce the third and fourth terms in Eq. (3). These terms align the unconditional outputs of the finetuned model with those of the original model across both stages.

3.3. The Heavy Hitters Among the Parameters

After determining the optimization objective, identifying the most effective parameters to optimize for achieving improved performance efficiently becomes crucial. Previous approaches typically divide the model into multiple modules, such as residual blocks, self-attention, or cross-attention, and select an entire module for finetuning [19, 20, 55, 99]. Among these, finetuning cross-attention modules is most common.

Inspired by saliency map techniques [13, 14, 25, 75, 76, 79], we propose a concept-specific saliency map enhanced by prompt and seed augmentation to precisely identify parameters suitable for finetuning. Compared to previous methods that compute the saliency map only once, we observe that the saliency map can vary depending on the prompt context and random seed, leading to instability. However, if we take the intersection of multiple saliency maps, the parameters within this intersection gradually become more stable and consistent as the number of maps increases (see Figure 4). This approach more accurately identifies the parameters responsible for the target concept, resulting in a consistent improvement in performance (as shown in the ablation study results in Table 5). Specifically, as illustrated in Figure 5, we first employ GPT-4 [61] to generate multiple prompts $\mathcal{C} = \{c_i\}_{i=1}^{N_c}$, each accompanied by a set of random seeds $\mathcal{S} = \{s_j\}_{j=1}^{N_s}$, to produce correspond-

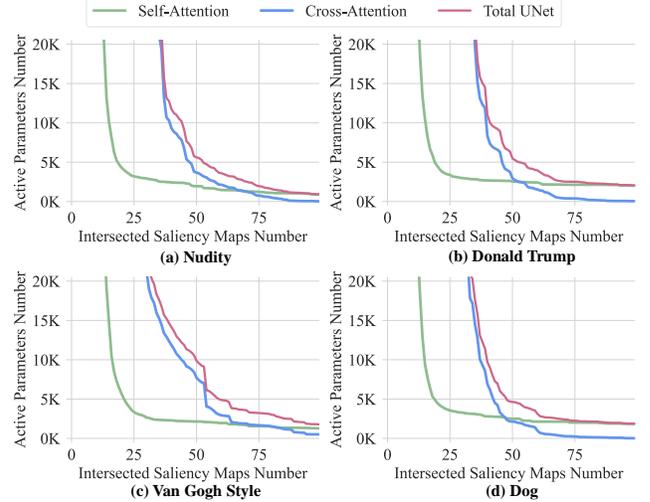


Figure 4. Each subplot shows the number of active parameters (y-axis) against the number of intersected saliency maps (x-axis) for four concepts: (a) Nudity, (b) Donald Trump, (c) Van Gogh Style, and (d) Dog. The number of active parameters converges across different concept types with around 100 intersected saliency maps.

ing gradient maps for the model parameters. By evaluating these gradients against a threshold, we obtain a set of weight saliency maps:

$$M_{c_i, s_j} = \mathbf{1}(|\nabla_{\theta} \mathcal{L}(z_{t_1}, z_{t_2}, t_1, t_2, c_i, s_j)| \geq \gamma), \quad (4)$$

where $\mathbf{1}(g \geq \gamma)$ is an element-wise indicator function that returns 1 for the i -th element if $g_i \geq \gamma$, and 0 otherwise; $|\cdot|$ denotes the element-wise absolute value operation; and $\gamma > 0$ is a predefined threshold. Each weight saliency map identifies critical parameters strongly correlated with the targeted concept across diverse prompt contexts. Finally, the intersection of these weight saliency maps obtained from various prompts and seeds yields the definitive concept-specific saliency map M^* :

$$M^* = \bigcap_{c_i \in \mathcal{C}} \bigcap_{s_j \in \mathcal{S}} M_{c_i, s_j}. \quad (5)$$

As a result, only a crucial subset of parameters is finetuned:

$$\theta \leftarrow \theta - \alpha \cdot M^* \odot \nabla_{\theta} \mathcal{L}(z_{t_1}, z_{t_2}, t_1, t_2, c_i, s_j), \quad (6)$$

where α is the learning rate and \odot denotes the element-wise multiplication. Intuitively, this mechanism identifies and finetunes only those parameters consistently influential for erasing the undesired concept across diverse conditions. Concept-specific saliency map M^* significantly narrows down the finetuning parameters, effectively preventing unnecessary perturbations to parameters unrelated to the targeted concept.

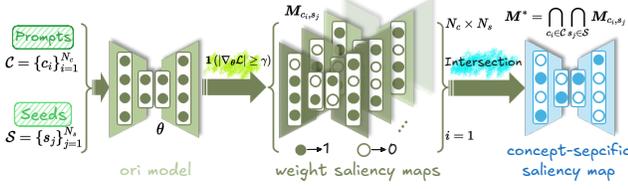


Figure 5. Generation of the concept-specific saliency map M^* . GPT-4 generates prompts $\mathcal{C} = \{c_i\}_{i=1}^{N_c}$, each paired with random seeds $\mathcal{S} = \{s_j\}_{j=1}^{N_s}$, which are used to compute gradient maps. After thresholding, saliency maps are obtained, and their intersection across all prompts and seeds yields M^* .

3.4. Boosting the Performance of Multi-Concept Erasure Frameworks

Our proposed trajectory-aware loss function seamlessly integrates with existing multi-concept erasure frameworks, such as MACE [55], offering a flexible and adaptable plug-and-play solution. Accordingly, it significantly boosts MACE’s performance in multi-concept scenarios, delivering new SOTA outcomes on tasks involving the erasure of 100 celebrity concepts and 100 artistic concepts.

As observed in MACE, erasing multiple concepts through either sequential or parallel finetuning often degrades performance. Sequential finetuning is susceptible to catastrophic forgetting, while parallel finetuning can lead to interference between concepts [55]. MACE addresses this by training a separate LoRA module for each concept to be erased, and subsequently fusing all LoRA modules into the cross-attention layers using a closed-form solution.

By integrating our loss function into the MACE framework, the initial training stage can be omitted. In the second stage, we replace MACE’s attention loss with our trajectory-aware loss to train individual LoRA modules ΔW_i for each concept, eliminating the need for the large Grounded-SAM model. After training all LoRA modules, we use the following objective function to fuse them into the cross-attention layers:

$$\min_{W^*} \sum_{i=1}^q \sum_{j=1}^p \left\| W^* \cdot e_j^f - (W + \Delta W_i) \cdot e_j^f \right\|_2^2 + \beta \sum_{j=p+1}^{p+m} \left\| W^* \cdot e_j^p - W \cdot e_j^p \right\|_2^2, \quad (7)$$

where W denotes the original weight matrix of either the key or value projection. The embedding e_j^f corresponds to concept-related tokens that we aim to erase, while e_j^p represents embeddings of unrelated, prior-preservation tokens. Here, q is the number of concepts to be erased, and p and m denote the numbers of targeted concept tokens and prior-preservation tokens, respectively.

As shown in Figure 6, the objective is to find a solu-

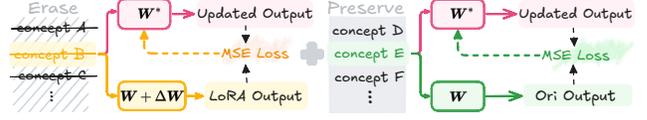


Figure 6. Multi-LoRA fusion for multi-concept erasure.

tion W^* that integrates multiple LoRA matrices, optimized for effective multi-concept erasure. This optimization problem has a closed-form solution [55]. Table 2 shows that our trajectory-aware loss function seamlessly integrates with the MACE framework for multi-concept erasure, substantially enhancing its performance.

4. Experiments

In this section, we present a comprehensive evaluation of our proposed method by benchmarking it against SOTA baselines on both single-concept erasure (NSFW removal; Section 4.2) and multi-concept erasure tasks, including 100-celebrity erasure (Section 4.3) and 100-artistic style erasure (Section 4.4). Finally, we perform ablation studies (Section 4.5) to assess the contribution of key components in our approach.

4.1. Implementation Details

We finetune all models based on Stable Diffusion (SD) v1.4 and generate outputs using the DDIM sampler [77] over 50 inference steps. Our experimental setup follows the settings described in MACE [55]. Each LoRA module undergoes 50 gradient update steps during training. For the baselines, we adopt the configurations provided in their respective original implementations.

4.2. Erasing NSFW Content

Configuration. In this experiment, we focus on removing the concept “nudity” from the model, representing a typical NSFW category. Specifically, we follow the “nudity”, “naked”, “erotic”, “sexual” prompts introduced in [28, 55] to guide the construction of the concept-specific saliency map M^* over the UNet. Based on M^* , we finetune SD v1.4 to eliminate the concept.

For evaluation, we use the full set of 4,703 prompts from the I2P dataset [73] along with their corresponding random seeds to generate images. We then apply NudeNet [65] with the threshold of 0.6 to detect exposed body parts in the sampled images, treating the detection results as an indicator of residual nudity in the model’s output. In addition, we assess the effectiveness of concept removal techniques in preserving benign content, utilizing the MS-COCO dataset [51]. We sample 30,000 captions from the validation split to generate images and compute FID [63] and CLIP score [66] as metrics for image quality and semantic alignment.

Table 1. Results of Erasing NSFW Content. The left side shows the number of exposed body parts detected on the I2P dataset using the NudeNet detector, while the right side presents the FID and CLIP on the COCO dataset. M: Male. F: Female.

Method	Inappropriate Image Prompt (I2P)									MS-COCO 30K	
	Armpits	Belly	Buttocks	Feet	Breasts (F)	Genitalia (F)	Breasts (M)	Genitalia (M)	Total ↓	FID ↓	CLIP ↑
FMN [99]	43	117	12	59	155	17	19	2	424	13.52	30.39
ESD-x [19]	59	73	12	39	100	6	18	8	315	14.41	30.69
ESD-u [19]	32	30	2	19	27	3	8	2	123	15.10	30.21
SLD-M [73]	47	72	3	21	39	1	26	3	212	16.34	30.90
AC [40]	153	180	45	66	298	22	67	7	838	14.13	31.37
SA [28]	72	77	19	25	83	16	0	0	292	-	-
EA [21]	-	-	-	-	-	-	-	-	199	21.75	30.24
UCE [20]	29	62	7	29	35	5	11	4	182	14.07	30.85
Receler[32]	39	26	5	10	13	1	12	9	115	-	-
MACE [55]	17	19	2	39	16	2	9	7	111	13.42	29.41
AdvUnlearn[100]	12	7	4	13	6	2	0	8	52	15.35	29.3
RealEra[52]	19	6	2	37	23	4	0	2	93	-	-
SPEED [49]	20	42	7	3	29	2	5	5	113	37.82	26.29
SalUn [13]	2	14	0	14	7	2	7	5	51	-	-
CE-SDWV [84]	13	46	2	2	13	0	1	6	84	13.66	30.80
SPM [57]	22	4	9	12	4	0	0	5	56	-	-
RECE [23]	17	23	0	8	8	0	6	4	66	-	-
SDD [37]	14	4	7	3	8	1	0	4	41	-	-
DuMo [24]	8	6	2	7	1	4	0	6	34	-	-
ACE [87]	5	7	3	6	2	3	4	9	39	14.69	30.80
Ours	1	5	2	4	8	2	0	1	23	14.44	30.64
SD v1.4	148	170	29	63	266	18	42	7	743	14.04	31.34
SD v2.1	105	159	17	60	177	9	57	2	586	14.87	31.53

Results Analysis. The experimental results are presented in Table 1. Our method generates significantly less NSFW content under the I2P benchmark prompts compared to other baselines, especially in challenging regions such as breasts. At the same time, our method also achieves competitive performance in terms of FID and CLIP scores. These results demonstrate that our approach can effectively remove explicit content from the model without compromising image quality.

4.3. Erasing Celebrity

Configuration. In this section, we evaluate the performance of our method on the task of simultaneously erasing multiple celebrity concepts, using the 200-celebrity dataset from MACE [55], which includes 100 celebrity concepts designated for erasure and 100 concepts intended to be preserved.

We conduct experiments by finetuning SD v1.4 to erase all 100 celebrity identities in the erasure group. We evaluate the effectiveness of our method by generating portraits of the targeted celebrities. Successful erasure is indicated by a low top-1 accuracy from GIPHY Celebrity Detector (GCD) [27] in identifying the erased identities. Additionally, to investigate the influence of our method on celebrities in the preservation group, we generate and evaluate their portraits in the same manner, where a high top-1 GCD accuracy reflects minimal impact on these preserved identities. We also report the harmonic mean H_c metric introduced

Table 2. Results of Erasing Celebrity. We report the accuracy for erased celebrities (Acc_e), accuracy for preserved celebrities (Acc_p), harmonic mean metric (H_c) and the proportion of clearly recognizable faces (Face Ratio). FID and CLIP are results based on MS-COCO dataset. SD v1.4 and SD v2.1 are used as reference base models.

Method	Acc_e ↓	Acc_p ↑	H_c ↑	Face Ratio↑	FID↓	CLIP↑
FMN [99]	0.9223	0.9076	0.1431	0.9940	13.95	31.31
ESD-x [19]	0.2784	0.2793	0.4027	0.8088	14.65	28.90
ESD-u [19]	0.0406	0.3909	0.4598	0.4724	15.14	29.02
SLD-M [73]	0.8706	0.7946	0.2237	0.9093	17.54	30.93
AC [40]	0.8913	0.9096	0.1977	0.9932	13.92	31.23
UCE [20]	0.0012	0.3790	0.5495	0.7179	106.57	19.17
RECE [23]	0.0243	0.2371	0.3816	-	177.57	12.09
SPEED [49]	0.0587	0.8554	0.8963	-	44.97	26.22
MACE [55]	0.0430	0.8456	0.8979	0.9820	12.82	30.21
Ours	0.0430	0.8807	0.9173	0.9816	11.71	30.40
SD v1.4	0.9648	0.9388	-	0.9876	14.04	31.34
SD v2.1	0.9324	0.9293	-	0.9879	14.87	31.53

in [55], which provides a balanced evaluation of the trade-off between successful erasure of unwanted celebrity concepts and the preservation of unrelated ones:

$$H_c = \frac{1}{(1 - Acc_e)^{-1} + (Acc_p)^{-1}}, \quad (8)$$

where H_c is the harmonic mean for celebrity erasure, Acc_e is the accuracy for the erased celebrities, and Acc_p for the preserved ones.

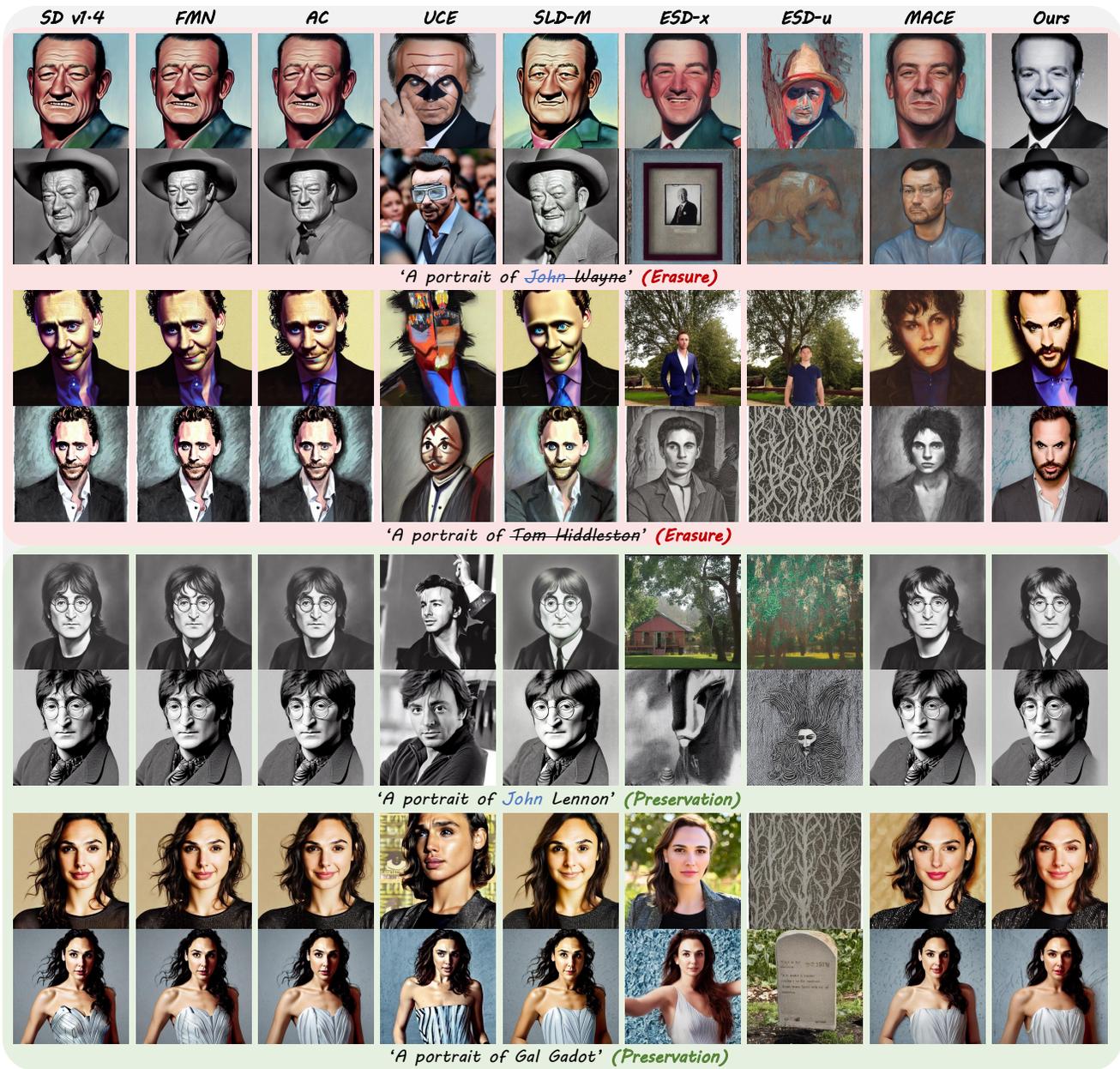


Figure 7. Qualitative comparison of erasing 100 celebrities from SD v1.4. John Wayne and Tom Hiddleston are in the erasure group for evaluating erasure performance; John Lennon and Gal Gadot are in preservation group for assessing preservation performance. Preserving John Lennon is challenging due to the shared first name with John Wayne.

Results Analysis. Figure 7 shows the qualitative comparison. Table 2 summarizes the performance of baselines on the celebrity concept erasure task. Our method achieves the highest H_c , outperforming all baselines and highlighting an excellent balance between concept erasure and preservation of unrelated ones.

Our method obtains the lowest FID score, surpassing all compared baselines and even the original SD models. A plausible reason for this improvement is that our finetun-

ing process, while primarily intended for erasing specific concepts, implicitly regularizes the model by encouraging more consistent representations of general concepts. Additionally, the CLIP score remains competitive, indicating minimal disruption to semantic alignment.

4.4. Erasing Art Style

Configuration. For art style erasure, we follow a similar training procedure as described in Section 4.3, with certain

Table 3. Results of Erasing 100 Art Styles. We report the CLIP score for erased artistic style ($CLIP_e$), CLIP score for preserved artistic style ($CLIP_p$), the overall score (H_a). FID and CLIP are results based on MS-COCO dataset.

Method	$CLIP_e \downarrow$	$CLIP_p \uparrow$	$H_a \uparrow$	FID-COCO \downarrow	CLIP-COCO \uparrow
FMN [99]	29.63	28.90	-0.73	13.99	31.31
ESD-x [19]	20.89	21.21	0.32	15.19	29.52
ESD-u [19]	19.66	19.55	-0.11	17.07	27.76
SLD-M [73]	28.49	27.89	-0.60	17.95	30.87
AC [40]	29.26	28.54	-0.72	14.08	31.29
UCE [20]	21.31	25.70	4.39	77.72	19.17
MACE [55]	22.59	28.58	5.99	12.71	29.51
Ours	20.6	26.78	6.18	12.96	27.63
SD v1.4	29.63	28.90	-	14.04	31.34

hyperparameter adjustments detailed in the Appendix. To evaluate performance, we use the 200-artist dataset from MACE [55], which consists of two groups: an erasure group of 100 artists whose styles are targeted for removal, and a preservation group of 100 artists whose styles are intended to be retained.

We use the CLIP score to assess how well the generated images align with the intended artistic style. For the erasure group, a lower CLIP score ($CLIP_e$) indicates better performance, as it suggests more effective removal of the target concept. In contrast, for the preservation group, a higher CLIP score ($CLIP_p$) is desirable, as it reflects minimal disruption to unrelated concepts. The overall performance is captured by $H_a = CLIP_p - CLIP_e$, where a higher value indicates better balance between preservation and erasure.

Results Analysis. Table 3 summarizes the performance of our method in erasing artistic styles. Our method achieves the highest H_a , substantially surpassing all baseline methods, demonstrating superior balance in effectively removing targeted art styles and preserving unrelated art styles. Considering the overall performance across other metrics, our strategy shows notable competitiveness compared to existing approaches.

4.5. Ablation Study

To investigate the contribution of key components in our approach, we conduct ablation studies on both multiple concepts (celebrity removal) and single concept (NSFW removal) tasks. The experimental configurations and corresponding results are presented in Tables 4 and 5, respectively.

We begin by ablating each component of our loss function in the context of celebrity removal. **Config A**, which applies $\mathcal{L}_{\text{erase}}$ across all stages without preserving the early-stage score function field, shows strong removal capability but clearly suffers in terms of preservation. **Config B** builds on Config A by adding $\mathcal{L}_{\text{uncond}}$ to maintain the unconditional score function, resulting in improved overall performance in

Table 4. Ablation study on multiple concepts (celebrity) removal. $\mathcal{L}_{\text{erase}}^*$: $\mathcal{L}_{\text{erase}}$ is applied at all denoising timesteps during training. $\mathcal{L}_{\text{erase}}$: $\mathcal{L}_{\text{erase}}$ is applied only during the mid-to-late stages of the denoising process in training.

Config	Components					Metrics		
	$\mathcal{L}_{\text{erase}}$	$\mathcal{L}_{\text{erase}}^*$	$\mathcal{L}_{\text{preserve}}$	$\mathcal{L}_{\text{uncond-early}}$	$\mathcal{L}_{\text{uncond-late}}$	$Acc_e \downarrow$	$Acc_p \uparrow$	$H_c \uparrow$
A	✗	✓	✗	✗	✗	0.0192	0.7785	0.8680
B	✗	✓	✗	✓	✓	0.0042	0.7848	0.8778
C	✓	✗	✗	✗	✗	0.0309	0.8094	0.8821
D	✓	✗	✗	✗	✓	0.0075	0.8013	0.8867
E	✓	✗	✓	✗	✗	0.0910	0.8545	0.8809
Ours		✗	✓	✓	✓	0.0430	0.8807	0.9173

Table 5. Ablation study on single concept (NSFW) removal. Single Map: M^* is generated using a single prompt and one random seed. Multi Maps: M^* is generated taking the intersection of saliency maps obtained using multiple prompts and multiple random seeds.

Config	Components		Inappropriate Image Prompt (I2P)			
	Single Map	Multi Maps	Breasts (M&F)	Genitalia (M&F)	Others	Total \downarrow
F	✗	✗	136	11	148	295
G	✓	✗	83	56	184	323
Ours	✗	✓	8	3	12	23

terms of H_c . Next, **Config C** applies $\mathcal{L}_{\text{erase}}$ only during the mid-to-late sampling stages, aiming to avoid disruption of the early-stage score function field. While this slightly improves H_c , the results remain unsatisfactory. **Config D** enhances Config C by adding $\mathcal{L}_{\text{uncond}}$, applied over the same timesteps as $\mathcal{L}_{\text{erase}}$, which further improves overall performance. **Config E** extends Config C by introducing $\mathcal{L}_{\text{preserve}}$, which helps retain the original score function field and significantly boosts preservation performance in terms of Acc_p . Our full method builds upon Config E by applying $\mathcal{L}_{\text{uncond}}$ across all stages, resulting in superior overall performance.

For NSFW content removal, **Config F** finetunes the entire UNet, while **Config G** finetunes only a subset of parameters using a saliency map obtained from a single calculation. However, Config G performs worse than Config F, suggesting that a saliency map generated from a single pass may be inaccurate. In contrast, our method derives a more precise concept-specific saliency map by taking the intersection of multiple saliency maps calculated from different prompts and seeds. This allows us to more accurately identify the parameters strongly associated with the concept, leading to substantially improved performance.

5. Conclusion

Our work introduces a geometric perspective on concept erasure within diffusion models. Utilizing this perspective, we found that reversing the condition direction of classifier-free guidance during the mid-to-late stages of the denoising process allows for modifying detailed content without

compromising the overall structural integrity of the generated images. Inspired by this insight, we propose ANT, a novel framework that effectively balances the removal of unwanted concepts while preserving unrelated elements. ANT demonstrates superior performance in both single- and multi-concept erasure scenarios, significantly outperforming current SOTA methods.

References

- [1] Stability AI. Stable diffusion v2.1 and dreamstudio updates 7-dec 22, 2022. [14](#)
- [2] Lucas Beerens, Alex D Richardson, Kaicheng Zhang, and Dongdong Chen. On the vulnerability of concept erasure in diffusion models. *arXiv preprint arXiv:2502.17537*, 2025. [1, 14](#)
- [3] Anh Bui, Long Vuong, Khanh Doan, Trung Le, Paul Montague, Tamas Abraham, and Dinh Phung. Erasing undesirable concepts in diffusion models with adversarial preservation. *arXiv preprint arXiv:2410.15618*, 2024. [1, 3, 14](#)
- [4] Anh Bui, Trang Vu, Long Vuong, Trung Le, Paul Montague, Tamas Abraham, Junae Kim, and Dinh Phung. Fantastic targets for concept erasure in diffusion models and where to find them. *arXiv preprint arXiv:2501.18950*, 2025. [1, 3, 14](#)
- [5] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. [1, 14](#)
- [6] Ruchika Chavhan, Da Li, and Timothy Hospedales. Conceptprune: Concept editing in diffusion models via skilled neuron pruning. *arXiv preprint arXiv:2405.19237*, 2024. [1, 3, 14](#)
- [7] Die Chen, Zhiwen Li, Mingyuan Fan, Cen Chen, Wenmeng Zhou, Yanhao Wang, and Yaliang Li. Growth inhibitors for suppressing inappropriate image concepts in diffusion models. In *The Thirteenth International Conference on Learning Representations*. [1](#)
- [8] Huiqiang Chen, Tianqing Zhu, Linlin Wang, Xin Yu, Longxiang Gao, and Wanlei Zhou. Safe and reliable diffusion models via subspace projection. *arXiv preprint arXiv:2503.16835*, 2025. [1, 3, 14](#)
- [9] Ruidong Chen, Honglin Guo, Lanjun Wang, Chenyu Zhang, Weizhi Nie, and An-An Liu. Trce: Towards reliable malicious concept erasure in text-to-image diffusion models. *arXiv preprint arXiv:2503.07389*, 2025. [1, 14](#)
- [10] Bartosz Cywiński and Kamil Deja. Saeuron: Interpretable concept unlearning in diffusion models with sparse autoencoders. *arXiv preprint arXiv:2501.18052*, 2025. [1, 14](#)
- [11] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. *Advances in Neural Information Processing Systems*, 35:16890–16902, 2022. [1, 14](#)
- [12] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. [14](#)
- [13] Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. *arXiv preprint arXiv:2310.12508*, 2023. [5, 7](#)
- [14] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018. [5](#)
- [15] Masane Fuchi and Tomohiro Takagi. Erasing concepts from text-to-image diffusion models with few-shot unlearning. *arXiv preprint arXiv:2405.07288*, 2, 2024. [1, 14](#)
- [16] Masane Fuchi and Tomohiro Takagi. Erasing with precision: Evaluating specific concept erasure from text-to-image generative models. *arXiv preprint arXiv:2502.13989*, 2025. [1, 3, 14](#)
- [17] Tatiana Gaintseva, Chengcheng Ma, Ziquan Liu, Martin Benning, Gregory Slabaugh, Jiankang Deng, and Ismail Elezi. Casteer: Steering diffusion models for controllable generation. *arXiv preprint arXiv:2503.09630*, 2025. [3, 14](#)
- [18] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. [14](#)
- [19] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. *arXiv preprint arXiv:2303.07345*, 2023. [1, 3, 4, 5, 7, 9](#)
- [20] 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*, pages 5111–5120, 2024. [1, 5, 7, 9, 14](#)
- [21] Daiheng Gao, Shilin Lu, Shaw Walters, Wenbo Zhou, Jiaming Chu, Jie Zhang, Bang Zhang, Mengxi Jia, Jian Zhao, Zhaoxin Fan, et al. Eraseanything: Enabling concept erasure in rectified flow transformers. *arXiv preprint arXiv:2412.20413*, 2024. [1, 3, 7, 14](#)
- [22] Hongcheng Gao, Tianyu Pang, Chao Du, Taihang Hu, Zhijie Deng, and Min Lin. Meta-unlearning on diffusion models: Preventing relearning unlearned concepts. *arXiv preprint arXiv:2410.12777*, 2024. [3, 14](#)
- [23] Chao Gong, Kai Chen, Zhipeng Wei, Jingjing Chen, and Yu-Gang Jiang. Reliable and efficient concept erasure of text-to-image diffusion models. In *European Conference on Computer Vision*, pages 73–88. Springer, 2024. [1, 7, 14](#)
- [24] Feng Han, Kai Chen, Chao Gong, Zhipeng Wei, Jingjing Chen, and Yu-Gang Jiang. Dumo: Dual encoder modulation network for precise concept erasure. *arXiv preprint arXiv:2501.01125*, 2025. [1, 3, 7, 14](#)
- [25] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149*, 2015. [5](#)

- [26] Tingxu Han, Weisong Sun, Yanrong Hu, Chunrong Fang, Yonglong Zhang, Shiqing Ma, Tao Zheng, Zhenyu Chen, and Zhenting Wang. Continuous concepts removal in text-to-image diffusion models. *arXiv preprint arXiv:2412.00580*, 2024. 1, 3, 14
- [27] Nick Hasty, Ihor Kroosh, Dmitry Voitek, and Dmytro Koruban. Giphy celebrity detector. <https://github.com/Giphy/celeb-detection-oss>. 7
- [28] Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep generative models. *arXiv preprint arXiv:2305.10120*, 2023. 1, 3, 6, 7, 14
- [29] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 2
- [30] Seunghoo Hong, Juhun Lee, and Simon S Woo. All but one: Surgical concept erasing with model preservation in text-to-image diffusion models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 21143–21151, 2024. 1
- [31] Yuepeng Hu, Zhengyuan Jiang, and Neil Zhenqiang Gong. Safetext: Safe text-to-image models via aligning the text encoder. *arXiv preprint arXiv:2502.20623*, 2025. 1, 14
- [32] Chi-Pin Huang, Kai-Po Chang, Chung-Ting Tsai, Yung-Hsuan Lai, Fu-En Yang, and Yu-Chiang Frank Wang. Receler: Reliable concept erasing of text-to-image diffusion models via lightweight erasers. In *European Conference on Computer Vision*, pages 360–376. Springer, 2024. 1, 3, 7, 14
- [33] Anubhav Jain, Yuya Kobayashi, Takashi Shibuya, Yuhta Takida, Nasir Memon, Julian Togelius, and Yuki Mitsufuji. Trasce: Trajectory steering for concept erasure. *arXiv preprint arXiv:2412.07658*, 2024. 1, 14
- [34] Abdullah Ayub Khan, Jing Yang, Asif Ali Laghari, Abdullah M Baqasah, Roobaea Alroobaea, Chin Soon Ku, Roohallah Alizadehsani, U Rajendra Acharya, and Lip Yee Por. Baiot-ems: Consortium network for small-medium enterprises management system with blockchain and augmented intelligence of things. *Engineering Applications of Artificial Intelligence*, 141:109838, 2025. 14
- [35] Changhoon Kim, Kyle Min, and Yezhou Yang. Race: Robust adversarial concept erasure for secure text-to-image diffusion model. In *European Conference on Computer Vision*, pages 461–478. Springer, 2024. 1, 14
- [36] Dahye Kim and Deepti Ghadiyaram. Concept steerers: Leveraging k-sparse autoencoders for controllable generations. *arXiv preprint arXiv:2501.19066*, 2025. 1, 14
- [37] Sanghyun Kim, Seohyeon Jung, Balhae Kim, Moonseok Choi, Jinwoo Shin, and Juho Lee. Towards safe self-distillation of internet-scale text-to-image diffusion models. *arXiv preprint arXiv:2307.05977*, 2023. 1, 3, 7, 14
- [38] Sanghyun Kim, Moonseok Choi, Jinwoo Shin, and Juho Lee. Safety alignment backfires: Preventing the re-emergence of suppressed concepts in fine-tuned text-to-image diffusion models. *arXiv preprint arXiv:2412.00357*, 2024.
- [39] Sanghyun Kim, Seohyeon Jung, Balhae Kim, Moonseok Choi, Jinwoo Shin, and Juho Lee. Safeguard text-to-image diffusion models with human feedback inversion. In *European Conference on Computer Vision*, pages 128–145. Springer, 2024. 1, 3, 14
- [40] 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*, pages 22691–22702, 2023. 1, 3, 7, 9
- [41] Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent space. *arXiv preprint arXiv:2210.10960*, 2022. 1
- [42] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. Accessed: February 21, 2025. 14
- [43] Byung Hyun Lee, Sungjin Lim, and Se Young Chun. Localized concept erasure for text-to-image diffusion models using training-free gated low-rank adaptation. *arXiv preprint arXiv:2503.12356*, 2025. 1, 14
- [44] Bin Li, Yixuan Weng, Bin Sun, and Shutao Li. A multi-tasking and multi-stage chinese minority pre-trained language model. In *China Conference on Machine Translation*, pages 93–105. Springer, 2022. 14
- [45] Bin Li, Yixuan Weng, Fei Xia, Bin Sun, and Shutao Li. Vpai.lab at medvidqa 2022: a two-stage cross-modal fusion method for medical instructional video classification. In *Proceedings of the 21st Workshop on Biomedical Language Processing*, pages 212–219, 2022.
- [46] Bin Li, Bin Sun, Shutao Li, Encheng Chen, Hongru Liu, Yixuan Weng, Yongping Bai, and Meiling Hu. Distinct but correct: generating diversified and entity-revised medical response. *Science China Information Sciences*, 67(3): 132106, 2024.
- [47] Bin Li, Yixuan Weng, Fei Xia, and Hanjun Deng. Towards better chinese-centric neural machine translation for low-resource languages. *Computer Speech & Language*, 84: 101566, 2024. 14
- [48] Feifei Li, Mi Zhang, Yiming Sun, and Min Yang. Detect-and-guide: Self-regulation of diffusion models for safe text-to-image generation via guideline token optimization. *arXiv preprint arXiv:2503.15197*, 2025. 1, 3, 14
- [49] Ouxiang Li, Yuan Wang, Xinting Hu, Houcheng Jiang, Tao Liang, Yanbin Hao, Guojun Ma, and Fuli Feng. Speed: Scalable, precise, and efficient concept erasure for diffusion models. *arXiv preprint arXiv:2503.07392*, 2025. 1, 7, 14
- [50] Shutao Li, Bin Li, Bin Sun, and Yixuan Weng. Towards visual-prompt temporal answer grounding in instructional video. *IEEE transactions on pattern analysis and machine intelligence*, 46(12):8836–8853, 2024. 14
- [51] 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*, pages 740–755. Springer, 2014. 6
- [52] Yufan Liu, Jinyang An, Wanqian Zhang, Ming Li, Dayan Wu, Jingzi Gu, Zheng Lin, and Weiping Wang. Realera: Semantic-level concept erasure via neighbor-concept mining. *arXiv preprint arXiv:2410.09140*, 2024. 1, 3, 7, 14

- [53] Zhili Liu, Kai Chen, Yifan Zhang, Jianhua Han, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan Yeung, and James T Kwok. Implicit concept removal of diffusion models. In *European Conference on Computer Vision*, pages 457–473. Springer, 2024. 1
- [54] Shilin Lu, Yanzhu Liu, and Adams Wai-Kin Kong. Tf-icon: Diffusion-based training-free cross-domain image composition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2294–2305, 2023. 14
- [55] 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*, pages 6430–6440, 2024. 1, 3, 5, 6, 7, 9, 14
- [56] Shilin Lu, Zihan Zhou, Jiayou Lu, Yuanzhi Zhu, and Adams Wai-Kin Kong. Robust watermarking using generative priors against image editing: From benchmarking to advances. *arXiv preprint arXiv:2410.18775*, 2024. 14
- [57] Mengyao Lyu, Yuhong Yang, Haiwen Hong, Hui Chen, Xuan Jin, Yuan He, Hui Xue, Jungong Han, and Guiguang Ding. One-dimensional adapter to rule them all: Concepts diffusion models and erasing applications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7559–7568, 2024. 1, 3, 7, 14
- [58] Zheling Meng, Bo Peng, Xiaochuan Jin, Yueming Lyu, Wei Wang, and Jing Dong. Concept corrector: Erase concepts on the fly for text-to-image diffusion models. *arXiv preprint arXiv:2502.16368*, 2025. 1, 3, 14
- [59] Quang H Nguyen, Hoang Phan, and Khoa D Doan. Unveiling concept attribution in diffusion models. *arXiv preprint arXiv:2412.02542*, 2024. 1, 14
- [60] 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. 1, 14
- [61] OpenAI. Hello gpt-4o, 2024. 5
- [62] Yong-Hyun Park, Sangdoon Yun, Jin-Hwa Kim, Junho Kim, Geonhui Jang, Yonghyun Jeong, Junghyo Jo, and Gayoung Lee. Direct unlearning optimization for robust and safe text-to-image models. *arXiv preprint arXiv:2407.21035*, 2024. 1, 14
- [63] 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*, pages 11410–11420, 2022. 6
- [64] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 14
- [65] platelminto. Nudenet: Neural nets for nudity classification, detection and selective censoring, 2023. 6
- [66] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 6
- [67] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 1, 14
- [68] Robin Rombach. Stable diffusion v1-4 model card. 2022. 14
- [69] 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*, pages 10684–10695, 2022. 1, 14
- [70] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023. 14
- [71] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 1, 14
- [72] Andrea Schioppa, Emiel Hoogeboom, and Jonathan Heek. Model integrity when unlearning with t2i diffusion models. *arXiv preprint arXiv:2411.02068*, 2024. 1, 3, 14
- [73] 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*, pages 22522–22531, 2023. 1, 3, 6, 7, 9
- [74] Reza Shirkavand, Peiran Yu, Shangqian Gao, Gowthami Somepalli, Tom Goldstein, and Heng Huang. Efficient fine-tuning and concept suppression for pruned diffusion models. *arXiv preprint arXiv:2412.15341*, 2024. 1, 3, 14
- [75] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013. 5
- [76] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017. 5
- [77] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 6
- [78] Koushik Srivatsan, Fahad Shamshad, Muzammal Naseer, and Karthik Nandakumar. Stereo: Towards adversarially robust concept erasing from text-to-image generation models. *arXiv preprint arXiv:2408.16807*, 2024. 1, 14
- [79] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3319–3328. JMLR. org, 2017. 5

- [80] Kartik Thakral, Tamar Glaser, Tal Hassner, Mayank Vatsa, and Richa Singh. Continual unlearning for foundational text-to-image models without generalization erosion. *arXiv preprint arXiv:2503.13769*, 2025. 1, 3, 14
- [81] Kartik Thakral, Tamar Glaser, Tal Hassner, Mayank Vatsa, and Richa Singh. Fine-grained erasure in text-to-image diffusion-based foundation models. *arXiv preprint arXiv:2503.19783*, 2025. 1, 14
- [82] Zhihua Tian, Sirun Nan, Ming Xu, Shengfang Zhai, Wenjie Qu, Jian Liu, Kui Ren, Ruoxi Jia, and Jiaheng Zhang. Sparse autoencoder as a zero-shot classifier for concept erasing in text-to-image diffusion models. *arXiv preprint arXiv:2503.09446*, 2025. 1, 3, 14
- [83] 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. 14
- [84] Jiahang Tu, Qian Feng, Chufan Chen, Jiahua Dong, Hanbin Zhao, Chao Zhang, and Hui Qian. Ce-sdwv: Effective and efficient concept erasure for text-to-image diffusion models via a semantic-driven word vocabulary. *arXiv preprint arXiv:2501.15562*, 2025. 1, 3, 7, 14
- [85] Ruipeng Wang, Junfeng Fang, Jiaqi Li, Hao Chen, Jie Shi, Kun Wang, and Xiang Wang. Ace: Concept editing in diffusion models without performance degradation. *arXiv preprint arXiv:2503.08116*, 2025. 1, 3, 14
- [86] Yanghao Wang and Long Chen. Improving diffusion-based data augmentation with inversion spherical interpolation. *arXiv preprint arXiv:2408.16266*, 2024. 1
- [87] Zihao Wang, Yuxiang Wei, Fan Li, Renjing Pei, Hang Xu, and Wangmeng Zuo. Ace: Anti-editing concept erasure in text-to-image models. *arXiv preprint arXiv:2501.01633*, 2025. 1, 3, 7, 14
- [88] Jing Wu and Mehrtash Harandi. Munba: Machine unlearning via nash bargaining. *arXiv preprint arXiv:2411.15537*, 2024. 1, 3, 14
- [89] Tianwei Xiong, Yue Wu, Enze Xie, Zhenguo Li, and Xihui Liu. Editing massive concepts in text-to-image diffusion models. *arXiv preprint arXiv:2403.13807*, 2024. 1, 14
- [90] Yuyang Xue, Edward Moroshko, Feng Chen, Steven McDonagh, and Sotirios A Tsaftaris. Crce: Coreference-retention concept erasure in text-to-image diffusion models. *arXiv preprint arXiv:2503.14232*, 2025. 1, 3, 14
- [91] Jing Yang, Liangyu Li, Lip Yee Por, Sami Bourouis, Sami Dhahbi, and Abdullah Ayub Khan. Harnessing multimodal data and deep learning for comprehensive gait analysis in pediatric cerebral palsy. *IEEE Transactions on Consumer Electronics*, 2024. 14
- [92] Jing Yang, Nika Anoosha Borojeni, Mehran Kazemi Charhardeh, Lip Yee Por, Roohallah Alizadehsani, and U Rajendra Acharya. A dual-method approach using autoencoders and transductive learning for remaining useful life estimation. *Engineering Applications of Artificial Intelligence*, 147:110285, 2025.
- [93] Jing Yang, Ke Tian, Huayu Zhao, Zheng Feng, Sami Bourouis, Sami Dhahbi, Abdullah Ayub Khan, Mouhebed-dine Berrima, and Lip Yee Por. Wastewater treatment monitoring: Fault detection in sensors using transductive learning and improved reinforcement learning. *Expert Systems with Applications*, 264:125805, 2025.
- [94] Jing Yang, Yuanguai Wu, Yuping Yuan, Haozhong Xue, Sami Bourouis, Mahmoud Abdel-Salam, Sunil Prajapat, and Lip Yee Por. Llm-ae-mp: Web attack detection using a large language model with autoencoder and multilayer perceptron. *Expert Systems with Applications*, 274:126982, 2025. 14
- [95] Tianyun Yang, Juan Cao, and Chang Xu. Pruning for robust concept erasing in diffusion models. *arXiv preprint arXiv:2405.16534*, 2024. 1, 3, 14
- [96] Xinlei Yu, Ahmed Elazab, Ruiquan Ge, Hui Jin, Xinchen Jiang, Gangyong Jia, Qing Wu, Qinglei Shi, and Changmiao Wang. Ich-scnet: Intracerebral hemorrhage segmentation and prognosis classification network using clip-guided sam mechanism. In *2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 2795–2800. IEEE, 2024. 14
- [97] Xinlei Yu, Xinyang Li, Ruiquan Ge, Shibin Wu, Ahmed Elazab, Jichao Zhu, Lingyan Zhang, Gangyong Jia, Taosheng Xu, Xiang Wan, et al. Ichpro: Intracerebral hemorrhage prognosis classification via joint-attention fusion-based 3d cross-modal network. In *2024 IEEE International Symposium on Biomedical Imaging (ISBI)*, pages 1–5. IEEE, 2024.
- [98] Xinlei Yu, Ahmed Elazab, Ruiquan Ge, Jichao Zhu, Lingyan Zhang, Gangyong Jia, Qing Wu, Xiang Wan, Lihua Li, and Changmiao Wang. Ich-prnet: a cross-modal intracerebral haemorrhage prognostic prediction method using joint-attention interaction mechanism. *Neural Networks*, 184:107096, 2025. 14
- [99] Eric Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not: Learning to forget in text-to-image diffusion models. *arXiv preprint arXiv:2303.17591*, 2023. 1, 3, 5, 7, 9
- [100] Yimeng Zhang, Xin Chen, Jinghan Jia, Yihua Zhang, Chongyu Fan, Jiancheng Liu, Mingyi Hong, Ke Ding, and Sijia Liu. Defensive unlearning with adversarial training for robust concept erasure in diffusion models. *Advances in Neural Information Processing Systems*, 37:36748–36776, 2024. 1, 7, 14
- [101] Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. In *European Conference on Computer Vision*, pages 385–403. Springer, 2024. 14
- [102] Mengnan Zhao, Lihe Zhang, Xingyi Yang, Tianhang Zheng, and Baocai Yin. Advanchor: Enhancing diffusion model unlearning with adversarial anchors. *arXiv preprint arXiv:2501.00054*, 2024. 1, 3, 14
- [103] Yuanzhi Zhu, Ruiqing Wang, Shilin Lu, Junnan Li, Hanshu Yan, and Kai Zhang. Ofsr: One-step flow for image super-resolution with tunable fidelity-realism trade-offs. *arXiv preprint arXiv:2412.09465*, 2024. 14

Appendix

A. Additional Related Work

With the advancement of deep learning [34, 44–47, 50, 91–94, 96–98] and generative models [5, 11, 54, 56, 60, 67, 69, 71, 103], an increasing number of studies have begun to focus on the issue of concept erasure in generative models.

A.1. Balancing Erasure and Preservation.

With the advancement of concept erasure techniques, the community has come to recognize that concept erasure should not only focus on the target concept but also aim to minimize the impact on unrelated concepts during fine-tuning. Numerous studies [3, 4, 6–8, 16, 17, 21, 22, 24, 26, 28, 32, 37, 39, 48, 52, 55, 57, 58, 72, 74, 80–82, 84, 85, 87, 88, 90, 95, 102] emphasize the model’s balanced performance between the target concept and unrelated concepts. MACE [55] introduces concept-focal importance sampling and modular LoRA integration, allowing for scalable multi-concept erasure while avoiding interference across modules. Several works [52, 90] explore semantic-aware preservation by modeling relationships between erased and retained concepts, improving quality retention in adjacent concept spaces. Some frameworks [81, 88] formalize the forgetting–retention trade-off, offering principled mechanisms to control degradation.

A.2. Finetuning Efficiency.

In addition to balancing erasure and preservation, several recent methods [3, 15, 20, 23, 36, 49, 57, 58, 85] have increasingly emphasized finetuning efficiency to meet practical demands. [20, 49] achieve erasure across hundreds of concepts within seconds by leveraging low-rank adapters or null-space constraints, enabling rapid adaptation across diffusion model variants. [23, 85] introduce closed-form or structure-aware updates that reduce erasure time by orders of magnitude. These advancements demonstrate a trend toward minimal-latency, high-throughput concept erasure that enables practical integration into production-scale text-to-image pipelines.

A.3. Scalability.

With the growing demand for safe and policy-compliant generative models, scalable multi-concept erasure techniques [8, 9, 20, 37, 43, 49, 55, 89] have emerged as a key direction in diffusion model editing. [20] introduces an editing framework that supports the simultaneous modification of multiple concepts through lightweight model updates. [55] leverages modular LoRA-based editing combined with closed-form integration to eliminate over 100 concepts with minimal interference. [89] adopts a two-stage process involving self-distillation and multi-layer editing, scaling up to 1,000 concepts while preserving

Table 6. Training hyperparameters for NSFW content, celebrity and art style erasure tasks.

Erasure Type	Learning Rate	Epochs	λ_1	λ_2	λ_3	t'
NSFW Content	5.0×10^{-4}	250	1.0	0.5	0.5	43
Celebrity	5.0×10^{-4}	400	0.4	0.5	0.2	40
Art Style	5.0×10^{-4}	400	0.4	0.5	0.2	47

specificity and visual fidelity. Additional methods such as [8, 9] enhance scalability through embedding-space operations or adversarially robust training objectives. These techniques collectively push concept erasure toward broader, more practical deployment scenarios requiring high-volume, reliable editing.

A.4. Robustness.

Despite successful concept suppression, erased models remain vulnerable to adversarial prompts that can reactivate undesirable content. A growing number of methods [2, 6, 9, 10, 17, 22, 23, 31–33, 35, 36, 43, 48, 58, 59, 62, 78, 82, 84, 87, 88, 95, 100] have begun to address this issue explicitly, aiming to improve model reliability in the face of prompt-based attacks. Methods such as [33, 78] tackle this by pairing adversarial prompt discovery with robust erasure objectives or inference-time steering, offering stronger defense without retraining. Others, like [35, 62, 100] incorporate adversarial training or preference-based optimization directly into the unlearning process to improve stability against attack. Complementary strategies from [10, 31, 59] focus on interpretable attribution or encoder-level alignment to neutralize unsafe inputs at their origin. Together, these works underscore the need for erasure techniques that are not only effective but resilient under adversarial conditions.

B. Hyperparameters Setup

Table 6 presents the specific hyperparameters used in the experiments for erasing different types of concepts.

C. Limitations and Future Work

Our work has primarily been tested on UNet-based diffusion models [1, 68]. As diffusion models increasingly adopt architectures like MMDiT [12, 42, 64], evaluating the compatibility of our approach with these new frameworks will be a key focus of our next phase. Additionally, assessing the robustness of our framework against adversarial prompts [83, 101] and its ability to withstand methods for learning personalized concepts [18, 70] will be of critical importance.

D. Additional Qualitative Results

Figure 8 presents a qualitative comparison of art style erasure and the preservation of unrelated concepts across dif-

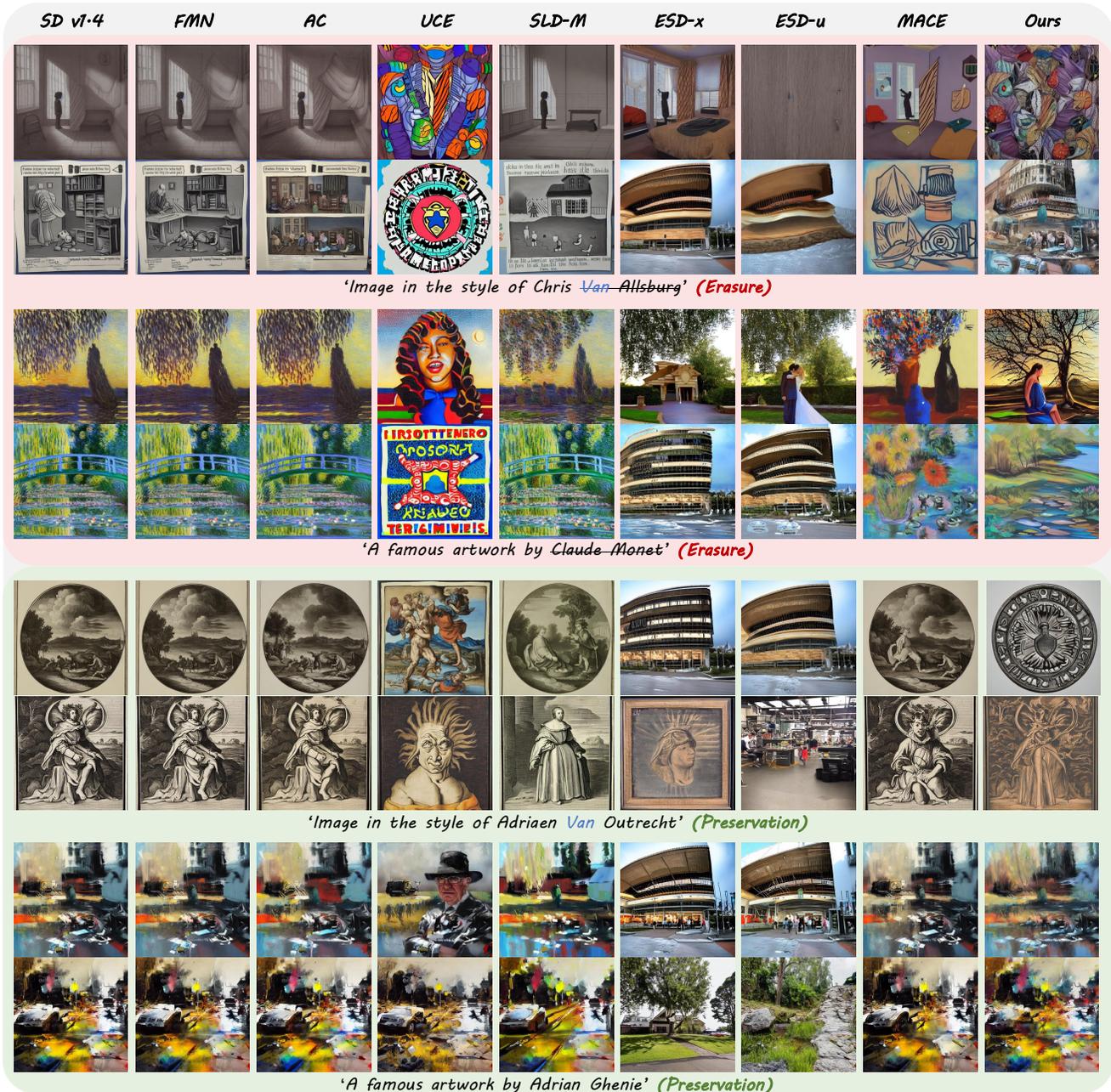


Figure 8. Qualitative comparison on art style erasure. The images on the same row are generated using the same random seed. Chris Van Allsburg and Claude Monet are in the erasure group, while Adriaen Van Utrecht and Adrian Ghenie are in the retention group.

ferent baselines. In the erasure rows, our approach effectively eliminates the target artistic styles (Chris Van Allsburg and Claude Monet) while retaining high-quality, plausible generation. In the preservation rows, our method successfully maintains the visual characteristics of unrelated artists (Adriaen Van Utrecht and Adrian Ghenie), showing minimal unintended impact on non-target styles.