

Separable Multi-Concept Erasure from Diffusion Models

Mengnan Zhao, Lihe Zhang, Tianhang Zheng, Yuqiu Kong and Baocai Yin

<https://github.com/Dlut-lab-zmn/SepCE4MU>

Abstract

Large-scale diffusion models, known for their impressive image generation capabilities, have raised concerns among researchers regarding social impacts, such as the imitation of copyrighted artistic styles. In response, existing approaches turn to machine unlearning techniques to eliminate unsafe concepts from pre-trained models. However, these methods compromise the generative performance and neglect the coupling among multi-concept erasures, as well as the concept restoration problem. To address these issues, we propose a Separable Multi-concept Eraser (SepME), which mainly includes two parts: the generation of concept-irrelevant representations and the weight decoupling. The former aims to avoid unlearning substantial information that is irrelevant to forgotten concepts. The latter separates optimizable model weights, making each weight increment correspond to a specific concept erasure without affecting generative performance on other concepts. Specifically, the weight increment for erasing a specified concept is formulated as a linear combination of solutions calculated based on other known undesirable concepts. Extensive experiments indicate the efficacy of our approach in eliminating concepts, preserving model performance, and offering flexibility in the erasure or recovery of various concepts.

1. Introduction

The field of text-to-image generation has witnessed remarkable development [52, 11, 24, 37], especially the occurrence of diffusion models (DMs) like DALL-E2 [36] and Stable Diffusion [38]. As the integration of DMs into practical applications [48, 20, 21] proves advantageous, addressing challenges related to their societal impact increasingly attracts the attention of researchers [3, 19, 12, 29]. One crucial challenge arises from diverse training data sources, potentially leading to unsafe image generation [50, 10], such as the creation of violent content or the imitation of specific artist styles. To resolve this concern, the machine unlearning (MU) technique has been proposed [51, 27, 45, 40],

which involves erasing the impact of specific data points or concepts to enhance model security, without necessitating complete retraining from scratch.

The recent MU research such as Erased Stable Diffusion (ESD) [10], Forget-me-not (FMN) [50], Safe self-distillation diffusion (SDD) [23], and Ablation Concept (AbConcept) [27], can be broadly categorized into untargeted concept erasure (e.g., FMN) and targeted concept erasure (e.g., ESD, SDD and AbConcept). Specifically, FMN minimizes the attention maps of forgotten concepts. In contrast, ESD, SDD, and AbConcept align the denoising distribution of forgotten concepts with a predefined distribution.

Despite recent advancements in MU [16, 33], there exist several drawbacks. Firstly, prior efforts concentrate on concept erasure, leading to considerable performance degradation in generative capability. Secondly, current erasure procedures are confined to single-concept elimination and pose challenges when extending them to multi-concept erasure. The multi-concept erasure can take two forms: simultaneous erasure of multiple concepts and iterative erasure of multiple concepts. The former means that multiple forgotten concepts are known in advance, while the latter implies that each erasure step only possesses knowledge of its previously forgotten concepts. Lastly, to the best of our knowledge, the concept restoration issue has not been considered. For instance, after the erasure of multiple artistic styles, the model owner may regain the copyright associated with some erased styles.

To address these issues, we propose an innovative Separable Multi-concept Eraser (SepME) that contains the generation of concept-irrelevant representations (G-CiRs) and the weight decoupling (WD). Specifically, G-CiRs aims to preserve overall model performance while effectively erasing undesirable concepts c_f through early stopping and weight regularization. Early stopping prevents the unlearning of substantial information irrelevant to c_f when $\Delta_\epsilon^{\theta_{\text{unlearn}}}$ and $\Delta_\epsilon^{\theta_{\text{ori}}}$ become irrelevant. Here, we define the difference of noise ϵ predicted by DMs with and without concept c_f as the representations of c_f , i.e., $\Delta_\epsilon^{\theta^*}$. θ_{ori} and

θ_{unlearn} are weights of the original and unlearned DMs, respectively. The regularization term restricts the deviation of θ_{unlearn} from θ_{ori} .

Considering the multi-concept erasure and subsequent restoration, WD characterizes the weight variation as $\Delta\theta_{1\sim N}$. Each independent weight increment $\Delta\theta_i$ is crafted to erase a specific concept $c_{i,f}$ without compromising the generation performance of models regarding other concepts. More precisely, the weight increment for erasing a specified concept is expressed as a linear combination of particular solutions calculated based on other known undesirable concepts. Each weight increment shares the same non-zero positions but has distinct values. These values are determined by the pre-calculated particular solutions and optimizable linear combination weights.

Our main contributions are summarized as follows: **(1)** To the best of our knowledge, the scenarios of multi-concept erasure and concept restoration have not been explored in previous literature. This work fills in these critical gaps and designs a separable multi-concept eraser; **(2)** To effectively unlearn undesirable concepts while maintaining overall model performance, our framework characterizes concept-irrelevant representations; **(3)** Through extensive experiments, we demonstrate that our method can improve erasing performance, maintain model generation capabilities, and offer flexibility in combining various forgotten concepts, encompassing both deletion and recovery.

2. Related Works

The image generation field has experienced rapid development in recent years, evolving from autoencoder [34, 31, 42], generative adversarial networks [6, 30, 28], unconditional diffusion models (DMs) [15, 5] to DMs enhanced with large-scale pre-trained image-text models [13, 46, 22] like CLIP [35]. These text-guided DMs, exemplified by DALL-E 2 [36] and Stable Diffusion [38], exhibit an excellent generative ability across various prompts c . The constraint for training DMs is formulated as

$$\mathcal{L}_{\text{DM}} = \mathbb{E}_{x_t \in \mathcal{D}, c, t, \epsilon_{\text{GT}} \in \mathcal{N}(0, \mathbf{I})} [\|\epsilon_{\text{GT}} - \epsilon_{\theta_{\text{dm}}}(x_t, c, t)\|_2^2],$$

where x_t represents the noised data or the noised latent representation [25], $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$. $\bar{\alpha}_t$ is the noise variance schedule. x_0 denotes the original reference image, and $\epsilon \in \mathcal{N}(0, \mathbf{I})$. \mathcal{D} is the training dataset. ϵ_{GT} means the ground truth noise. $\epsilon_{\theta_{\text{dm}}}(x_t, c, t)$ denotes the t -th step noise predicted by DMs. $\|\cdot\|_2^2$ is the squared ℓ_2 -norm function. Additionally, researchers have indicated the unknown concept generation capability of DMs through the fine-tuning of partial model weights on small reference sets [39, 18, 47, 9]. Nevertheless, DMs also induce potential risks associated with privacy violations and copyright infringement, such as the training data leakage [1, 7, 8], the

imitation of various artistic styles [44, 41], and the generation of sensitive content [49]. Consequently, there is a growing focus on erasing specific outputs from pre-trained DMs [3, 27].

Existing research primarily falls into three distinct directions: removal of unsafe data and model retraining [4], integration of additional plug-ins to guide model outputs [2, 32], and fine-tuning of model weights through MU techniques [10, 50, 27]. The drawback of the first direction is that large-scale model retraining demands considerable computational resource and time. The risk of the second direction is that, with the public availability of model structures and weights, malicious users can easily remove plug-ins. This work focuses on the third direction, *i.e.*, machine unlearning.

The majority of unlearning methods for DMs can be summarized as:

$$\epsilon_{\theta_{\text{op}}}(x_t, c_f, t) \leftarrow \begin{cases} \epsilon_{\text{target}} & \text{if } x_0 \in \mathcal{D}_f \\ \epsilon_{\text{GT}} & \text{otherwise,} \end{cases} \quad (1)$$

where θ_{op} represents optimizable model weights, *e.g.*, the parameters of cross-attention modules in DMs. x_t can be obtained through either the diffusion process or the sampling process. $\epsilon_{\theta_{\text{op}}}(x_t, c_f, t)$ denotes the noise predicted by unlearned DMs at the t -th step. \mathcal{D}_f refers to the dataset containing the forgotten concept c_f . ϵ_{target} and ϵ_{GT} represent the noise of predefined target concepts and the ground-truth noise added in the diffusion process, respectively.

For instance, ESD [10] leverages the predicted noise for both concept-free c_\emptyset and forgotten concepts c_f to construct ϵ_{target} ,

$$\epsilon_{\text{target}} = (1 + \eta)\epsilon_{\theta_{\text{dm}}}(x_t, c_\emptyset, t) - \eta\epsilon_{\theta_{\text{dm}}}(x_t, c_f, t),$$

where θ_{dm} represents parameters of the frozen DMs. η is the hyperparameter. SDD [23] directly maps the prediction distribution of erased concepts c_f to the prediction distribution of concept-free c_\emptyset , $\epsilon_{\text{target}} = \epsilon_{\theta_{\text{dm}}}(x_t, c_\emptyset, t)$. AbConcept [27] assigns anchor concepts c^* for each erased concept c_f , such as c_f is ‘Van Gogh painting’ and c^* is ‘painting’ or c_f is ‘a photo of Grumpy cat’ and c^* is ‘a photo of cat’, $\epsilon_{\text{target}} = \epsilon_{\theta_{\text{dm}}}(x_t, c^*, t)$. In contrast, FMN [50] is an untargeted concept erasure method, which minimizes attention weights corresponding to the forgotten concepts c_f .

These advanced approaches focus on unlearning concepts but compromise model performance significantly. Moreover, they have not considered the scenarios of multi-concept erasure and subsequent restoration. In this work, we introduce a separable multi-concept erasure framework. It incorporates an untargeted concept-irrelevant erasure mechanism to preserve model performance during concept erasure, and a weight decoupling mechanism to provide flexibility in both the erasure and recovery of concepts.

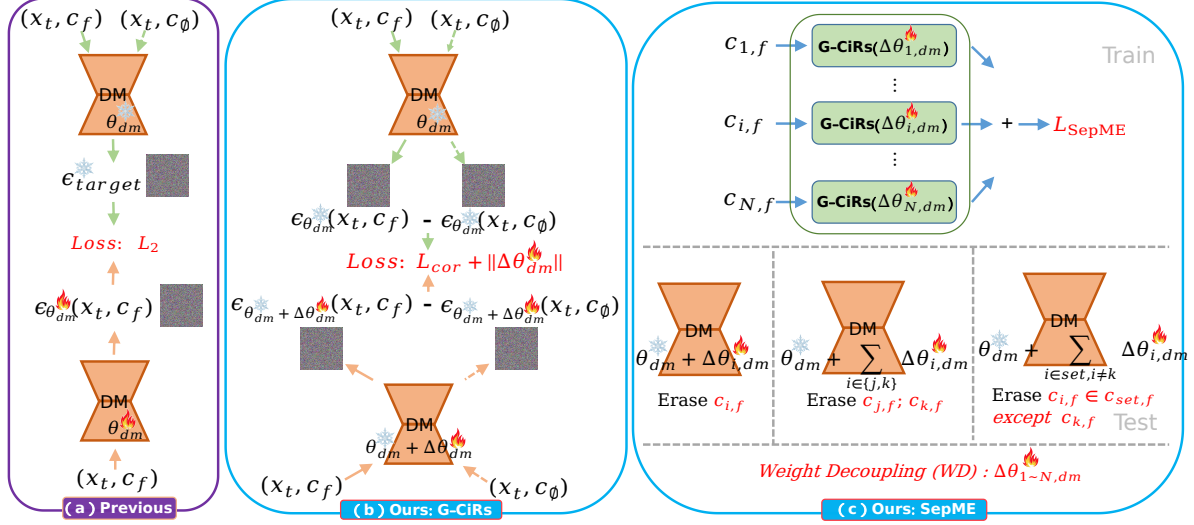


Figure 1. Overview of various unlearning techniques for DMs. ‘Ice flowers’ and ‘flames’ represent frozen and optimizable model weights respectively. ϵ_{target} is the predefined noise distribution. \mathcal{L}_2 denotes the ℓ_2 -norm function. x_t means noised samples. c_f and c_0 signify the forgotten and blank prompts, respectively. \mathcal{L}_{cor} is the correlation function. SepME separates optimizable weights as $\Delta\theta_{1\sim N, \text{dm}}$. Each $\Delta\theta_{i, \text{dm}}$ is designed to erase a specific concept $c_{i, f}$ without compromising the generation performance of models regarding other concepts.

3. Proposed Method

Our proposed separable multi-concept eraser (SepME) aims to flexibly erase or recover multiple concepts while preserving the overall model performance. An overview of SepME is illustrated in Fig. 1, which incorporates the generation of concept-irrelevant representations (G-CiRs) and the weight decoupling (WD).

3.1. G-CiRs

To maintain the generative capability of diffusion models (DMs) for regular concepts (or concept-free c_0) during unlearning, G-CiRs prevents the erasure of significant but irrelevant information to forgotten concepts $c_{i, f} \in c_f$, $i \in [1, N]$, where N is the number of erased concepts. Specifically, given original DMs parameterized by θ_{dm} , we employ the noise difference $\Delta_\epsilon(c_{i, f}, \theta_{\text{dm}}) = \epsilon_{\theta_{\text{dm}}}(x_t, c_{i, f}, t) - \epsilon_{\theta_{\text{dm}}}(x_t, c_0, t)$ to represent the concept $c_{i, f}$. To successfully erase concepts c_f from DMs, the representations of concepts c_f for unlearned and original DMs should be uncorrelated, that is $\forall_{i \in [1, N]} \mathcal{L}_{\text{cor}}(c_{i, f}, \Delta\theta_{\text{dm}}) = 0$,

$$\mathcal{L}_{\text{cor}}(c_{i, f}, \Delta\theta_{\text{dm}}) = \text{Avg}(\Delta_\epsilon(c_{i, f}, \theta_{\text{dm}}) \odot \Delta_\epsilon(c_{i, f}, \theta_{\text{dm}} + \Delta\theta_{\text{dm}})), \quad (2)$$

where $\Delta\theta_{\text{dm}}$ represents the learnable weight increments of unlearned DMs, \odot denotes the element-wise product, and $\text{Avg}(\cdot)$ calculates the average value. Eq. (2) actually computes the relevance between two representations for the concept $c_{i, f}$. On this basis, we fine-tune $\Delta\theta_{\text{dm}}$ by

$$\min_{\Delta\theta_{\text{dm}}} \mathcal{L}_{\text{G-CiRs}} = \sum_{i=1}^N \eta_i \mathcal{L}_{\text{cor}}(c_{i, f}, \Delta\theta_{\text{dm}}) + \lambda \|\Delta\theta_{\text{dm}}\|_p, \quad (3)$$

$$\text{s.t. } \forall_{i \in [1, N]} \mathcal{L}_{\text{cor}}(c_{i, f}, \Delta\theta_{\text{dm}}) = 0,$$

where η_i is used to balance the losses of multiple concepts, $\eta_i = \frac{\|\mathcal{L}_{\text{cor}}(c_{1, f}, \Delta\theta_{\text{dm}})\|_2}{\|\mathcal{L}_{\text{cor}}(c_{i, f}, \Delta\theta_{\text{dm}})\|_2}$. λ denotes the hyperparameter. $\|\Delta\theta_{\text{dm}}\|_p$ restricts the weight deviation of the unlearned DMs from the original ones.

To satisfy the condition of zero relevance in Eq. (3), we utilize the momentum statistic method since the values of $\mathcal{L}_{\text{cor}}(c_{i, f}, \Delta\theta_{\text{dm}})$ computed from various noised samples x_t exhibit significant variations. Specifically, early stopping is activated once $\mathcal{L}_{\text{mom}}^n \leq \tau$, where τ denotes the threshold, with a default value of 0.

$$\mathcal{L}_{\text{mom}}^n = \alpha \mathcal{L}_{\text{mom}}^{n-1} + (1 - \alpha) \sum_{i=1}^N \eta_i \mathcal{L}_{\text{cor}}(c_{i, f}, \Delta\theta_{\text{dm}}), \quad (4)$$

where α is the hyper-parameter and n is the unlearning step.

3.2. Weight Decoupling (WD)

To resolve the subsequent restoration issue of multi-concept erasure, we decompose the weights $\Delta\theta_{\text{dm}}$ in Eq. (2) into $\Delta\theta_{1\sim N, \text{dm}}$ for flexibly manipulating various concepts. Each independent weight increment $\Delta\theta_{i, \text{dm}}$ aims to unlearn a specific concept $c_{i, f}$ without compromising the generation performance of models regarding other concepts. The process of separable erasure can be formulated as:

$$\min_{\Delta\theta_{1\sim N, \text{dm}}} \mathcal{L}_{\text{SepME}} = \sum_{i=1}^N \eta_i \mathcal{L}_{\text{cor}}(c_{i, f}, \Delta\theta_{i, \text{dm}}) + \lambda \|\Delta\theta_{1\sim N, \text{dm}}\|_p,$$

$$\text{s.t. } \forall_{i \in [1, N]} \text{cond}(c_{i, f}, \Delta\theta_{i, \text{dm}}), \quad (5)$$

where the conditions include:

$$\begin{aligned}
\mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{i,\text{dm}}) &== 0, \\
\epsilon_{\theta_{\text{dm}}}(x_t, c_{\emptyset}, t) &== \epsilon_{\theta_{\text{dm}} + \Delta\theta_{i,\text{dm}}}(x_t, c_{\emptyset}, t), \\
\forall_{j \in [1, N], j \neq i} \epsilon_{\theta_{\text{dm}}}(x_t, c_{j,f}, t) &== \epsilon_{\theta_{\text{dm}} + \Delta\theta_{i,\text{dm}}}(x_t, c_{j,f}, t), \\
\forall_{c_{\not\in f}} \epsilon_{\theta_{\text{dm}}}(x_t, c_{\not\in f}, t) &\approx \epsilon_{\theta_{\text{dm}} + \Delta\theta_{i,\text{dm}}}(x_t, c_{\not\in f}, t),
\end{aligned} \tag{6}$$

where x_t can be an arbitrary image and $c_{\not\in f}$ represents concepts that do not belong to c_f . The first condition aims to erase forgotten concepts, while the rest mitigate the impact of fine-tuning on other concepts.

To meet these conditions, we first need to identify the nonzero positions for $\Delta\theta_{i,\text{dm}}$. It is evident from Eq. (6) that these positions are image-independent. Consequently, only the to_k and to_v layers of the cross-attention modules are selected, which are exclusively designed for extracting text embeddings in DMs. Other positions, such as the to_q layer for extracting image embeddings and the FFN (Feed-Forward Network) for updating fused embeddings, have been fixed. For $\forall_{i \in [1, N]} \Delta\theta_{i,\text{dm}}$, they share the same nonzero positions but have distinct values.

Next, we take $\Delta\theta_{\text{to}_k}$ as an example to analyze how to determine its value, where $\forall \Delta\theta_{\text{to}_k} \in \Delta\theta_{i,\text{dm}}$. Notably, $\epsilon_{\theta_{\text{dm}}}(x_t, c_k, t) == \epsilon_{\theta_{\text{dm}} + \Delta\theta_{\text{to}_k}}(x_t, c_k, t)$ means $c_k \otimes \Delta\theta_{\text{to}_k} == 0$, where $c_k \in \mathbb{R}^{d_{\text{emb}} \times d_{\text{in}}}$, $\Delta\theta_{\text{to}_k} \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}}$. d_{emb} , d_{in} and d_{out} indicate feature dimensions. \otimes denotes matrix multiplication. Thus, Eq. (6) can be rewritten as

$$\begin{aligned}
\mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{i,\text{dm}}) &== 0, \\
c_{\emptyset} \otimes \Delta\theta_{\text{to}_k} &== 0, \\
\forall_{j \in [1, N], j \neq i} c_{j,f} \otimes \Delta\theta_{\text{to}_k} &== 0, \\
\forall_{c_{\not\in f}} c_{\not\in f} \otimes \Delta\theta_{\text{to}_k} \approx 0^{d_{\text{emb}} \times d_{\text{out}}} &\Rightarrow \Delta\theta_{\text{to}_k} \approx 0^{d_{\text{in}} \times d_{\text{out}}}.
\end{aligned} \tag{7}$$

◆ To make $\Delta\theta_{\text{to}_k}$ satisfy the second and third conditions in Eq. (7), we first compute the particular solutions S_p to a system of linear equations $A \otimes S_p = 0$,

$$A = [c_{\emptyset}^{\top}; c_{1,f}^{\top}; \dots; c_{i-1,f}^{\top}; c_{i+1,f}^{\top}; \dots; c_{N,f}^{\top}]^{\top}, \tag{8}$$

where A is a constant matrix for specified concepts c_f , $A \in \mathbb{R}^{(N \cdot d_{\text{emb}}) \times d_{\text{in}}}$. \top means matrix transpose. $S_p \in \mathbb{R}^{d_{\text{in}} \times d_{\text{in}} - r}$, where r is the rank of A , with $r \leq N \cdot d_{\text{emb}}$. $d_{\text{in}} - r$ quantifies the number of solutions within S_p . $d_{\text{in}} \gg d_{\text{emb}}$ in DMs. Notably, to remove the original biases in solutions S_p , we normalize each element of S_p to a unit vector.

Then, each column of $\Delta\theta_{\text{to}_k}$ can be formulated as a linear combination of these solutions S_p ,

$$\Delta\theta_{\text{to}_k} = (w \otimes S_p^{\top})^{\top}, \tag{9}$$

where w is an optimizable variable and represents the linear combination weights, $w \in \mathbb{R}^{d_{\text{out}} \times (d_{\text{in}} - r)}$.

◆ To further make $\Delta\theta_{\text{to}_k}$ satisfy the fourth condition in Eq. (7), we introduce a scaling factor β as follows,

$$\Delta\theta_{\text{to}_k} = (w \otimes (\beta S_p^{\top}))^{\top}. \tag{10}$$

Meanwhile, w is initialized to a zero matrix. Additionally, we replace $\|\Delta\theta_{1 \sim N, \text{dm}}\|_p$ in Eq. (5) with the $\|\mathcal{W}\|_p$. Here, \mathcal{W} represents the set of optimizable variables w defined for both to_k and to_v layers.

Overall, the objective of SepME is simplified as:

$$\begin{aligned}
\min_{\mathcal{W}} \mathcal{L}_{\text{SepME}} &= \sum_{i=1}^N \eta_i \mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{i,\text{dm}}) + \lambda \|\mathcal{W}\|_p, \\
s.t. \quad \forall_{i \in [1, N]} \mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{i,\text{dm}}) &== 0.
\end{aligned} \tag{11}$$

Evaluation for SepME. We combine various $\Delta\theta_{i,\text{dm}}$ to erase the corresponding concepts. For instance, DMs with $\theta_{\text{dm}} + \sum_{i \in \{j, k\}} \Delta\theta_{i,\text{dm}}$ eliminate the concepts $c_{j,f}$ and $c_{k,f}$.

4. Experiments

4.1. Experimental Settings

Implementation Details. We follow prior works [10, 23] to fine-tune Stable Diffusion [38]. The optimization process utilizes the Adam optimizer for a maximum of 1000 iterations and our early stopping strategy. The batch size is set equal to the number of erased concepts. When exclusively evaluating the G-CiRs module, the learning rate is set to 1e-6, and we opt to fine-tune the cross-attention modules. For assessing the SepME, the learning rate is adjusted to 1e-2, and optimization is conducted on the to_k and to_v layers of the cross-attention modules. The default values for hyperparameters τ , α in Eq. (4), β in Eq. (10), and λ in Eq. (11) are set to 0, 0.9, 1e-4, and 3e-5, respectively. The threshold τ controls the moment of early stopping, and ablation studies for τ are provided in the appendix. All experiments are executed on 2 RTX 3090 GPUs.

Evaluation metrics. The evaluation metrics include modifications to model parameters $\|\Delta\theta_{\text{dm}}\|_p = \frac{\|\Delta\theta_{\text{dm}}\|_1}{M}$, perceptual distance measured by Perceptual Image Patch Similarity (LPIPS), and classification accuracy (ACC). Here, M denotes the number of layers. LPIPS quantifies the similarity between the original image and the image generated by unlearned DMs, calculated based on AlexNet [26] with settings from the source code¹.

We calculate ACC using various pre-trained classification models. For the style classification model, we consider a blank concept and nine artist styles: ‘Van Gogh’, ‘Picasso’, ‘Cezanne’, ‘Jackson Pollock’, ‘Caravaggio’, ‘Keith Haring’, ‘Kelly McKernan’, ‘Tyler Edlin’, and ‘Kilian Eng’. In each category, we generate 1000 images using the original DMs with artist names (or ‘’) as prompts. 70% of the data is allocated for training purposes, while the remaining 30% is reserved for testing. Only the fully connected (FC) layer of the pre-trained ResNet18 model [14] is optimized with 20 epochs. The cyclical learning rate [43] is

¹<https://github.com/richzhang/PerceptualSimilarity>

Table 1. Quantitative results of the single concept erasure. ‘VG’, ‘PC’ and ‘CE’ are artists of ‘Van Gogh’, ‘Picasso’ and ‘Cezanne’, respectively. $\bar{\uparrow}$: Since the unlearning aims to erase the concept style, an intermediate value may indicate better performance. i in FMN_i represents the iteration step. Text in **red** and **blue** denotes the best and second-best results, respectively.

ACC/LPIPS		Unlearning methods					
Erased	Evaluated	ORI	FMN ₁₀	FMN ₂₀	ESD	AbConcept	Ours-G-CiRs
VG	VG*($\downarrow/\bar{\uparrow}$)	1.000/0.000	0.544/0.433	0.500/0.517	0.320/0.472	0.000 /0.469	0.228 /0.363
	PC \dagger (\uparrow/\downarrow)	1.000/0.000	1.000 /0.186	1.000 /0.256	1.000 /0.301	0.900/0.190	1.000 /0.175
	CE \dagger (\uparrow/\downarrow)	1.000/0.000	1.000 /0.178	0.964/0.257	1.000 /0.164	0.820/0.217	1.000 /0.154
PC	VG \dagger (\uparrow/\downarrow)	1.000/0.000	1.000 /0.205	0.952 /0.265	0.908/0.209	0.684/0.227	0.908/ 0.180
	PC*($\downarrow/\bar{\uparrow}$)	1.000/0.000	1.000/0.231	0.952/0.339	0.400/0.424	0.000 /0.397	0.052 /0.358
	CE \dagger (\uparrow/\downarrow)	1.000/0.000	1.000 /0.149	1.000 /0.175	1.000 /0.211	0.752/0.194	1.000 /0.189
CE	VG \dagger (\uparrow/\downarrow)	1.000/0.000	0.952 /0.201	0.772/0.305	0.908 /0.238	0.728/0.257	0.908 /0.276
	PC \dagger (\uparrow/\downarrow)	1.000/0.000	1.000 /0.184	1.000 /0.207	1.000 /0.314	0.852/0.251	1.000 /0.204
	CE*($\downarrow/\bar{\uparrow}$)	1.000/0.000	1.000/0.189	0.820/0.351	0.428/0.372	0.036 /0.360	0.000 /0.363
$\sum \cdot^* - \sum \cdot \dagger (\downarrow/\bar{\uparrow})$		-	-3.41/-0.25	-3.42/-0.26	-4.67/-0.17	-4.70 /-0.11	-5.54 /-0.09
$\ \Delta\theta_{dm}\ _p \downarrow$		ORI	FMN ₁₀	FMN ₂₀	ESD	AbConcept	Ours-G-CiRs
VG		0.000	44.18	86.05	120.8	158.3	43.17
PC		0.000	45.41	93.63	126.9	146.5	36.90
CE		0.000	45.31	92.32	128.5	156.9	36.55

Table 2. Quantitative results of the multi-concept erasure. ‘VG’, ‘PC’ and ‘CE’ are ‘Van Gogh’, ‘Picasso’ and ‘Cezanne’, respectively. $\bar{\uparrow}$: Since the unlearning aims to erase the concept style, an intermediate value may indicate better performance. i in FMN_i represents the iteration step. Text in **red** and **blue** denotes the best and second-best results, respectively.

ACC/LPIPS		Unlearning methods						
Erased	Evaluated	ORI	FMN ₂₀	FMN ₃₀	FMN ₅₀	ESD	AbConcept	Ours-G-CiRs
VG+PC+CE	VG*($\downarrow/\bar{\uparrow}$)	1./0.	0.544/0.487	0.500/0.511	0.136 /0.555	0.364/0.467	0.092 /0.421	0.224/0.426
	PC*($\downarrow/\bar{\uparrow}$)	1./0.	0.952/0.299	0.552/0.299	0.000 /0.436	0.252/0.442	0.132 /0.329	0.000 /0.359
	CE*($\downarrow/\bar{\uparrow}$)	1./0.	0.180 /0.359	0.000 /0.424	0.000 /0.505	0.356/0.343	0.500/0.286	0.000 /0.419
	Others \dagger (\uparrow/\downarrow)	1./0.	0.955/0.228	0.878/0.269	0.693/0.358	0.897/0.252	0.977 /0.198	0.958 /0.223
$\sum \cdot^* - \sum \cdot \dagger (\downarrow/\bar{\uparrow})$		-	0.721/0.917	-0.17/0.965	-0.56 /1.138	-0.07/1.000	-0.26/0.838	-0.73 /0.981
$\ \Delta\theta_{dm}\ _p \downarrow$		ORI	FMN ₂₀	FMN ₃₀	FMN ₅₀	ESD	AbConcept	Ours-G-CiRs
VG+PC+CE		0.0	92.58	150.3	243.5	128.4	153.8	58.22

employed with the maximum learning rate of 0.01. As for the object classification network, we directly utilize the pre-trained ResNet50.

Advanced unlearning methods, including FMN [50], ESD [10] and AbConcept [23], are used as **baselines**.

4.2. Style Removal

To evaluate the performance of our methods in eliminating styles, we focus on three prominent artists: ‘Van Gogh’, ‘Picasso’, and ‘Cezanne’. During fine-tuning, we employ the images for obtaining the style classification model as inference images x_0 to yield x_t , *i.e.*, $x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon$. During evaluation, we generate 250 images for each concept, *i.e.*, 50 seeds for each concept, and 5 images per seed. As our SepME relies on the proposed G-CiRs module, we first indicate the performance of G-CiRs and then validate the efficacy of SepME.

4.2.1 Evaluation for G-CiRs

For the *single concept erasure*, one style is chosen as the forgotten concept, while others serve as evaluation concepts. The quantitative results are presented in Tab. 1. **On one hand, our G-CiRs achieves optimal performance in terms of ACC and LPIPS metrics across three artistic styles. On the other hand, the proposed G-CiRs induces fewer modifications to model weights.** Furthermore, qualitative comparisons in Figs. 2~4 also demonstrate the efficacy of our method in erasing artistic style and preserving generative performance across other concepts. These visual samples are produced with artist names or sentences² containing artist names as prompts for various unlearned DMs.

Likewise, we compare our G-CiRs with previous works under the *simultaneous erasure of multiple concepts*.

²https://github.com/rohitgandikota/erasing/blob/main/data/art_prompts.csv

Table 3. Quantitative results of SepME when simultaneously fine-tuning $\Delta\theta_{1\sim 3, dm}$. ‘VG’, ‘CE’ and ‘PC’ are artists of ‘Van Gogh’, ‘Cezanne’ and ‘Picasso’, respectively. $\Delta\theta_{1, dm}$, $\Delta\theta_{2, dm}$, $\Delta\theta_{3, dm}$ are optimizable weights for erasing ‘VG’, ‘CE’ and ‘PC’, respectively. i in FMN_i represents the iteration step. For each combination, we employ AbConcept and G-CiRs as baselines to re-finetune all layers of the cross-attention module in DMs to eliminate the corresponding concepts. Text in **red** indicates the best result.

		SepME(ACC/LPIPS)							
θ_{dm}	+0	$+\Delta\theta_{1, dm}$	$+\Delta\theta_{2, dm}$	$+\Delta\theta_{3, dm}$	$+\sum_{i=1}^2 \Delta\theta_{i, dm}$	$+\sum_{i \in \{1, 3\}} \Delta\theta_{i, dm}$	$+\sum_{i=2}^3 \Delta\theta_{i, dm}$	$+\sum_{i=1}^3 \Delta\theta_{i, dm}$	
VG	1./0.	0.320/0.371	0.956/0.182	0.956/0.182	0.364/0.372	0.272/0.371	0.956/0.182	0.364/0.372	
CE	1./0.	1.000/0.144	0.180/0.303	1.000/0.144	0.180/0.303	1.000/0.144	0.180/0.303	0.180/0.303	
PC	1./0.	1.000/0.185	1.000/0.185	0.000/0.440	1.000/0.185	0.000/0.440	0.000/0.440	0.000/0.440	
		AbConcept(ACC/LPIPS)							
c_f	-	VG	CE	PC	VG+CE	VG+PC	CE+PC	VG+CE+PC	
VG	1./0.	0.000/0.469	0.728/0.257	0.684/0.227	0.136/0.402	0.044/0.430	0.728/0.257	0.000/0.425	
CE	1./0.	0.820/0.217	0.036/0.360	0.752/0.194	0.572/0.253	0.820/0.194	0.252/0.298	0.464/0.269	
PC	1./0.	0.900/0.190	0.852/0.251	0.000/0.397	0.952/0.172	0.352/0.284	0.728/0.257	0.152/0.301	
		G-CiRs(ACC/LPIPS)							
c_f	-	VG	CE	PC	VG+CE	VG+PC	CE+PC	VG+CE+PC	
VG	1./0.	0.228/0.363	0.908/0.180	0.908/0.276	0.184/0.384	0.184/0.382	0.700/0.259	0.224/0.426	
CE	1./0.	1.000/0.154	0.000/0.363	1.000/0.189	0.108/0.354	0.780/0.173	0.036/0.427	0.000/0.359	
PC	1./0.	1.000/0.175	1.000/0.204	0.052/0.358	1.000/0.210	0.152/0.295	0.000/0.449	0.000/0.419	

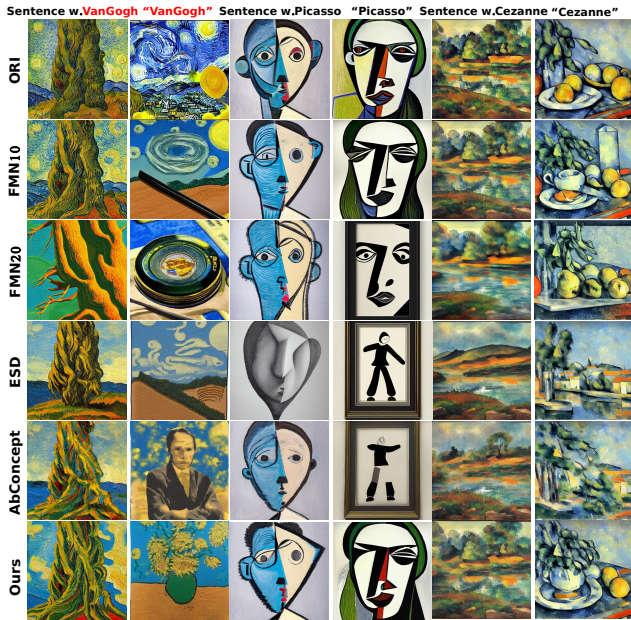


Figure 2. Qualitative comparison among various unlearning techniques for DMs with ‘Van Gogh’ as the erased concept.

Specifically, we consider the styles ‘Van Gogh,’ ‘Picasso,’ and ‘Cezanne’ as forgotten concepts and include additional styles in Sec. 4.1 for evaluation purposes. The experimental results in Tab. 2 demonstrate that our method achieves optimal performance when simultaneously erasing multiple concepts. Furthermore, we provide visual examples in Fig. 5. As observed, even when using forgotten concepts as prompts, our G-CiRs does not generate images containing erased styles. Notably, these generated images contain few discernible objects, which is expected given that the pro-

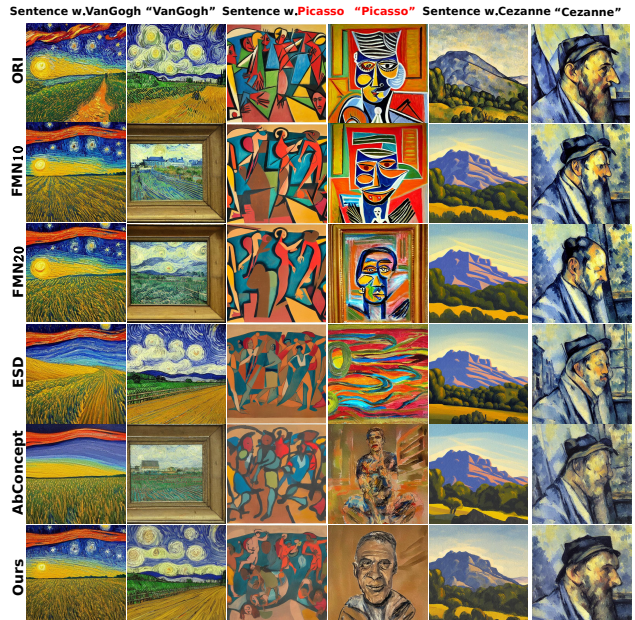


Figure 3. Qualitative comparison among various unlearning techniques for DMs with ‘Picasso’ as the erased concept.

posed G-CiRs is an untargeted unlearning technique.

4.2.2 Evaluation for SepME

The preceding experiments assessed the unlearning performance of various approaches. In the following, we investigate the concept restoration issue overlooked by these methods under two practical scenarios: unlearning multiple concepts simultaneously or iteratively unlearning multiple concepts. 1) The former knows all forgotten concepts for each concept erasure. 2) The latter only has knowledge of previ-

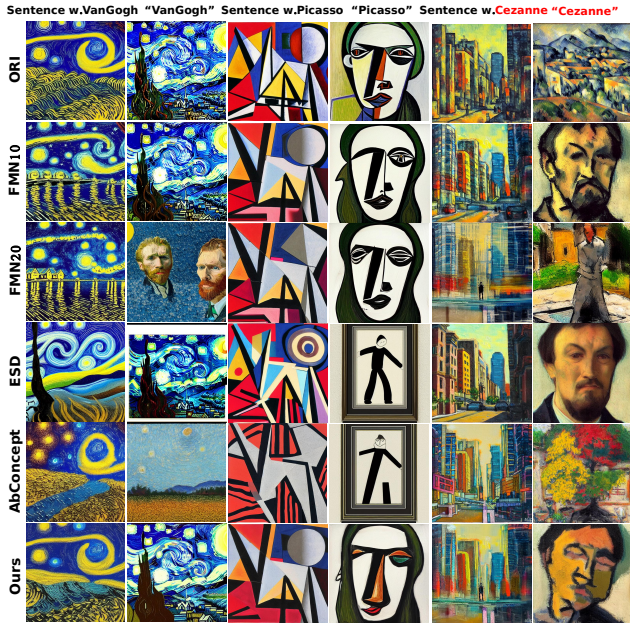


Figure 4. Qualitative comparison among various unlearning techniques for DMs with ‘Cezanne’ as the erased concept.

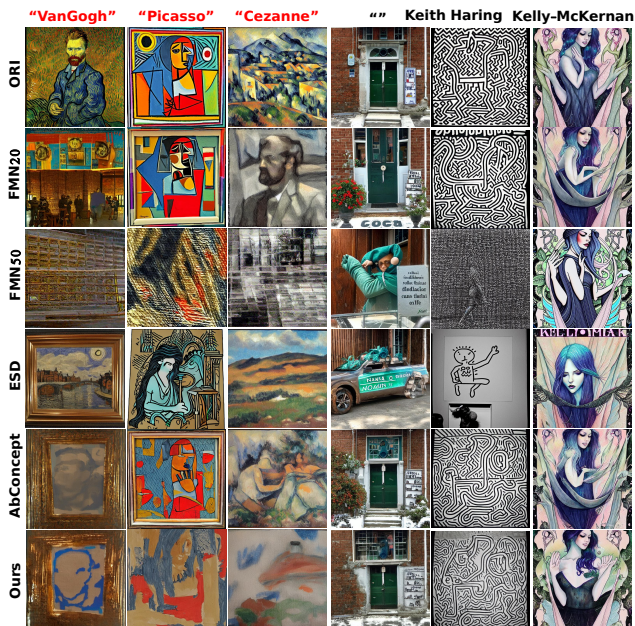


Figure 5. Qualitative comparison among various unlearning techniques under multi-concept erasure. ‘Ours’ indicates G-CiRs.

ously forgotten concepts at each erasure step.

After unlearning, various weights are randomly combined to erase corresponding concepts, such as $\theta_{dm} + \Delta\theta_{1,dm}$, $\theta_{dm} + \Delta\theta_{2,dm}$, and $\theta_{dm} + \Delta\theta_{3,dm}$ for erasing ‘Van Gogh,’ ‘Cezanne,’ and ‘Picasso,’ respectively. The combination $\theta_{dm} + \sum_{i=2}^3 \Delta\theta_{i,dm}$ erases the concepts ‘Cezanne’ and ‘Picasso’. As only the to_k and to_v layers of the



Figure 6. Visual examples of the proposed SepME.



Figure 7. Failure cases of the proposed SepME.

cross-attention modules are chosen as nonzero positions for SepME, achieving $\mathcal{L}_{cor}(\text{‘Cezanne’}, \Delta\theta_{2,dm}) = 0$ is challenging. Therefore, the threshold τ for ‘Cezanne’ is adjusted to $1.5e-4$.

Simultaneous erasure of multiple concepts. We first optimize $\Delta\theta_{1\sim3,dm}$ simultaneously. The quantitative and qualitative results are presented in Tab. 3 and Fig. 6, respectively. For each combination, we utilize AbConcept and G-CiRs as baselines to re-finetune all layers of the cross-attention module in DMs to eliminate the corresponding concepts. It can be observed that SepME can effectively and flexibly erase various concepts, achieving comparable performance to separately fine-tuned methods, such as the results of $\Delta\theta_{3,dm}$ and $\sum_{i=1}^3 \Delta\theta_{i,dm}$. Fig. 7 displays several failure cases, i.e., images produced by DMs with $\sum_{i=1}^3 \Delta\theta_{i,dm}$ but classified as forgotten concept categories. Additionally, the results of SepME on $\Delta\theta_{i,dm}$ and $\sum_{i \in \{j,k\}} \Delta\theta_{i,dm}$ indicate the feasibility of concept restoration after multi-concept erasures.

Next, we individually fine-tune $\Delta\theta_{1,dm}$, $\Delta\theta_{2,dm}$, and $\Delta\theta_{3,dm}$ before combining them to simultaneously eliminate

Table 4. Quantitative results of SepME when separately fine-tuning $\Delta\theta_{1\sim 3, \text{dm}}$. The concepts $c_{1,f}$, $c_{2,f}$, and $c_{3,f}$ correspond to ‘Van Gogh’, ‘Cezanne’, and ‘Picasso’, respectively. $\Delta\theta_{1, \text{dm}}$, $\Delta\theta_{2, \text{dm}}$, $\Delta\theta_{3, \text{dm}}$ are optimizable weights for erasing ‘VG’, ‘CE’ and ‘PC’, respectively. In SepME₁, $\Delta\theta_{1\sim 3, \text{dm}}$ are separately optimized when all forgotten concepts are known. SepME₂ follows the mode of iterative concept erasure, *i.e.*, the t -th erasure step only possesses knowledge (Kn) of the previously forgotten concepts.

SepME ₁ -(ACC/LPIPS)- $Kn = [c_{1,f}; c_{2,f}; c_{3,f}]$							
θ_{dm}	$+\Delta\theta_{1, \text{dm}}$	$+\Delta\theta_{2, \text{dm}}$	$+\Delta\theta_{3, \text{dm}}$	$+\sum_{i=1}^2 \Delta\theta_{i, \text{dm}}$	$+\sum_{i \in \{1,3\}} \Delta\theta_{i, \text{dm}}$	$+\sum_{i=2}^3 \Delta\theta_{i, \text{dm}}$	$+\sum_{i=1}^3 \Delta\theta_{i, \text{dm}}$
VG	0.356/0.364	1.000/0.182	0.908/0.182	0.308/0.365	0.308/0.365	<u>0.956/0.181</u>	0.308/0.364
CE	1.000/0.144	0.320/0.304	1.000/0.144	0.320/0.304	<u>1.000/0.144</u>	0.288/0.304	0.320/0.304
PC	1.000/0.185	1.000/0.185	0.000/0.460	<u>1.000/0.185</u>	0.000/0.460	0.000/0.460	0.000/0.460
SepME ₂ -(ACC/LPIPS)							
Kn	$[c_{1,f}]$	$[c_{1,f}; c_{2,f}]$	$[c_{1,f}; c_{2,f}; c_{3,f}]$	-	-	-	-
VG	0.228/0.363	0.956/0.182	0.956/0.182	0.228/0.363	0.228/0.363	<u>0.956/0.182</u>	0.228/0.363
CE	1.000/0.182	0.000/0.404	1.000/0.172	0.000/0.376	<u>1.000/0.177</u>	0.000/0.411	0.052/0.394
PC	0.964/0.150	1.000/0.144	0.288/0.277	<u>0.964/0.150</u>	0.108/0.292	0.288/0.277	0.108/0.292
Abconcept-(ACC/LPIPS)							
VG	0.000/0.469	0.728/0.257	0.684/0.227	0.000/0.472	0.000/0.469	<u>0.636/0.277</u>	0.000/0.487
CE	0.820/0.217	0.036/0.360	0.752/0.194	0.288/0.299	<u>0.680/0.234</u>	0.252/0.288	0.144/0.349
PC	0.900/0.190	0.852/0.251	0.000/0.397	<u>0.600/0.277</u>	0.152/0.304	0.200/0.334	0.052/0.364
G-CiRs-(ACC/LPIPS)							
VG	0.228/0.363	0.908/0.180	0.908/0.276	0.184/0.386	0.092/0.382	<u>0.544/0.276</u>	0.184/0.410
CE	1.000/0.154	0.000/0.363	1.000/0.189	0.108/0.261	<u>0.716/0.178</u>	0.216/0.313	0.036/0.331
PC	1.000/0.175	1.000/0.204	0.052/0.358	<u>1.000/0.220</u>	0.452/0.247	0.200/0.280	0.200/0.303

multiple concepts. The experimental results are detailed in Tab. 4. It is apparent that all unlearning methods effectively erase forgotten concepts. However, both AbConcept and our G-CiRs exhibit shortcomings in restoring forgotten concepts. For example, Abconcept and G-CiRs under the setting $\theta_{\text{dm}} + \sum_{i \in \{1,3\}} \Delta\theta_{i, \text{dm}}$ only perform 68% and 71.6% classification accuracy for ‘Van Gogh’, respectively. In contrast, our SepME₁ in Tab. 4 demonstrates nearly perfect recovery of erased concepts. [This emphasizes the effectiveness of our SepME in restoring forgotten concepts and the feasibility of separately optimizing \$\Delta\theta_{i, \text{dm}}\$.](#)

Iterative-concept erasure. To realize multi-concept erasure and concept restoration under this scenario, we avoid sequentially fine-tuning model weights, as restoring the early weights $\Delta\theta_{i, \text{dm}}$ inevitably affects the erasure performance of $\Delta\theta_{> i, \text{dm}}$. Inspired by the success of previous experiments where we optimized $\Delta\theta_{1\sim 3, \text{dm}}$ individually, we achieve iterative concept erasure through this optimization mode.

In the initial unlearning step, G-CiRs is employed to fine-tune all parameters of cross-attention modules in DMs. This enables better unlearning of the forgotten concept with smaller weight modifications. In each subsequent unlearning step t , we recalculate S_p using $c_{< t, f}$ to construct the weight increments of to_k and to_v layers. These increments are further fine-tuned to erase $c_{t, f}$. The experimental results presented as SepME₂ in Tab. 4 show comparable performance to SepME₁. [Overall, SepME can effectively achieve iterative concept erasure and concept restoration.](#)

4.3. Object Removal.

The experimental results on object removal yield similar conclusions to those on style removal. For detailed information, please refer to the appendix.

5. Conclusion

In this study, we present an innovative machine unlearning technique for diffusion models, namely separable multi-concept eraser (SepME). SepME leverages a correlation term and momentum statistics to yield concept-irrelevant representations. It not only maintains overall model performance during concept erasure but also adeptly balances loss magnitudes across multiple concepts. Furthermore, SepME allows for the separation of weight increments, providing flexibility in manipulating various concepts, including concept restoration and iterative concept erasure. Extensive experiments validate the effectiveness of our methods.

Broader Impact. As the field of deep learning continues to evolve, it presents both exciting opportunities and profound responsibilities for our community. While recent advances hold promise for solving complex problems, they also raise concerns regarding ethical and societal implications. As researchers in this domain, we recognize our obligation to comprehend and address the challenges associated with the widespread adoption of deep learning technology. Machine learning models, despite their potential benefits, can harbor harmful biases, unintended behaviors, and pose risks to user privacy. Our work contributes to this discourse by proposing a post-processing ‘unlearning’ phase aimed

at mitigating these concerns. Through extensive empirical investigation, we demonstrate progress over previous solutions in practical settings. However, it’s important to acknowledge that while our approach, SpeME, represents a significant step forward, we cannot claim perfect mitigation of these issues. Therefore, it’s imperative that caution is exercised in the practical application of deep learning techniques, and that rigorous auditing and evaluation of machine learning models are conducted.

References

- [1] Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 5253–5270, 2023.
- [2] 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. *arXiv preprint arXiv:2309.06135*, 2023.
- [3] Rishav Chourasia and Neil Shah. Forget unlearning: Towards true data-deletion in machine learning. In *International Conference on Machine Learning*, pages 6028–6073. PMLR, 2023.
- [4] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- [6] Xin Ding, Yongwei Wang, Zuheng Xu, William J Welch, and Z Jane Wang. Ccgan: Continuous conditional generative adversarial networks for image generation. In *International conference on learning representations*, 2020.
- [7] Jinhao Duan, Fei Kong, Shiqi Wang, Xiaoshuang Shi, and Kaidi Xu. Are diffusion models vulnerable to membership inference attacks? *arXiv preprint arXiv:2302.01316*, 2023.
- [8] Jan Dubiński, Antoni Kowalczyk, Stanisław Pawlak, Przemysław Rokita, Tomasz Trzcinski, and Paweł Morawiecki. Towards more realistic membership inference attacks on large diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4860–4869, 2024.
- [9] 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.
- [10] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
- [11] Songwei Ge, Taesung Park, Jun-Yan Zhu, and Jia-Bin Huang. Expressive text-to-image generation with rich text. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7545–7556, 2023.
- [12] Aditya Golatkar, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. Training data protection with compositional diffusion models. *arXiv preprint arXiv:2308.01937*, 2023.
- [13] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10696–10706, 2022.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in neural information processing systems*, pages 6840–6851, 2020.
- [16] Seunghoo Hong, Juhun Lee, and Simon S Woo. All but one: Surgical concept erasing with model preservation in text-to-image diffusion models. *arXiv preprint arXiv:2312.12807*, 2023.
- [17] Jeremy Howard and Sylvain Gugger. Fastai: A layered api for deep learning. *Information*, 11(2):108, 2020.
- [18] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [19] Chi-Pin Huang, Kai-Po Chang, Chung-Ting Tsai, Yung-Hsuan Lai, and Yu-Chiang Frank Wang. Receler: Reliable concept erasing of text-to-image diffusion models via lightweight erasers. *arXiv preprint arXiv:2311.17717*, 2023.
- [20] Bahjat Kavar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Con-*

- ference on Computer Vision and Pattern Recognition*, pages 6007–6017, 2023.
- [21] Amirhossein Kazerooni, Ehsan Khodapanah Aghdam, Moein Heidari, Reza Azad, Mohsen Fayyaz, Ilker Hacihaliloglu, and Dorit Merhof. Diffusion models for medical image analysis: A comprehensive survey. *arXiv preprint arXiv:2211.07804*, 2022.
- [22] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2426–2435, 2022.
- [23] 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.
- [24] Yunji Kim, Jiyoung Lee, Jin-Hwa Kim, Jung-Woo Ha, and Jun-Yan Zhu. Dense text-to-image generation with attention modulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7701–7711, 2023.
- [25] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [27] 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*, 2023.
- [28] Bowen Li, Xiaojuan Qi, Philip Torr, and Thomas Lukasiewicz. Lightweight generative adversarial networks for text-guided image manipulation. *Advances in Neural Information Processing Systems*, 33:22020–22031, 2020.
- [29] Hang Li, Chengzhi Shen, Philip Torr, Volker Tresp, and Jindong Gu. Self-discovering interpretable diffusion latent directions for responsible text-to-image generation. *arXiv preprint arXiv:2311.17216*, 2023.
- [30] Rui Liu, Yixiao Ge, Ching Lam Choi, Xiaogang Wang, and Hongsheng Li. Divco: Diverse conditional image synthesis via contrastive generative adversarial network. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16377–16386, 2021.
- [31] Zhi-Song Liu, Wan-Chi Siu, and Li-Wen Wang. Variational autoencoder for reference based image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 516–525, 2021.
- [32] Ninareh Mehrabi, Palash Goyal, Christophe Dupuy, Qian Hu, Shalini Ghosh, Richard Zemel, Kai-Wei Chang, Aram Galstyan, and Rahul Gupta. Flirt: Feedback loop in-context red teaming. *arXiv preprint arXiv:2308.04265*, 2023.
- [33] Zixuan Ni, Longhui Wei, Jiacheng Li, Siliang Tang, Yueting Zhuang, and Qi Tian. Degeneration-tuning: Using scrambled grid shield unwanted concepts from stable diffusion. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 8900–8909, 2023.
- [34] Taesung Park, Jun-Yan Zhu, Oliver Wang, Jingwan Lu, Eli Shechtman, Alexei Efros, and Richard Zhang. Swapping autoencoder for deep image manipulation. *Advances in Neural Information Processing Systems*, 33:7198–7211, 2020.
- [35] 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.
- [36] 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.
- [37] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [38] 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.
- [39] 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.
- [40] 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.

- [41] Shawn Shan, Jenna Cryan, Emily Wenger, Haitao Zheng, Rana Hanocka, and Ben Y Zhao. Glaze: Protecting artists from style mimicry by text-to-image models. *arXiv preprint arXiv:2302.04222*, 2023.
- [42] Huajie Shao, Shuochao Yao, Dachun Sun, Aston Zhang, Shengzhong Liu, Dongxin Liu, Jun Wang, and Tarek Abdelzaher. Controlvae: Controllable variational autoencoder. In *International Conference on Machine Learning*, pages 8655–8664. PMLR, 2020.
- [43] Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE winter conference on applications of computer vision (WACV)*, pages 464–472. IEEE, 2017.
- [44] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6048–6058, 2023.
- [45] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [46] Xingqian Xu, Zhangyang Wang, Gong Zhang, Kai Wang, and Humphrey Shi. Versatile diffusion: Text, images and variations all in one diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7754–7765, 2023.
- [47] Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhensu Chen, Xiaopeng Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language models. *arXiv preprint arXiv:2309.14717*, 2023.
- [48] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4):1–39, 2023.
- [49] Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Nan Xu, and Qiang Xu. Mma-diffusion: Multimodal attack on diffusion models. *arXiv preprint arXiv:2311.17516*, 2023.
- [50] 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.
- [51] 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. *arXiv preprint arXiv:2310.11868*, 2023.
- [52] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. Towards language-free training for text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17907–17917, 2022.

A. Code availability

The code is available at <https://github.com/Dlut-lab-zmn/SepCE4MU>.

B. Algorithms

The algorithmic details for generating concept-irrelevant representations are outlined in Alg. 1. Furthermore, Alg. 2 provides a comprehensive explanation of the separable multi-concept eraser.

Algorithm 1: G-CiRs.

Input: The diffuser $G(\cdot)$, the frozen weights θ_{dm} , the weight increment $\Delta\theta_{\text{dm}}$, the N forgotten concepts $c_{i,f} \in c_f$, the blank prompt c_\emptyset , the inference dataset $x_0 \in D$, the noise schedule $\bar{\alpha}_t$, the hyperparameter λ .

Output: The fine-tuned model increment $\Delta\theta_{\text{dm}}$.

```
1 for  $n, x_0 \in D$  do
2   Randomly select a sampling step  $t$ ;
3    $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \epsilon \in \mathcal{N}(0, \mathbf{I})$ ;
4    $\epsilon_{c_f} = G_{\theta_{\text{dm}}}(x_t, c_f, t)$ ;
5    $\epsilon_{c_\emptyset} = G_{\theta_{\text{dm}}}(x_t, c_\emptyset, t)$ ;
6    $\epsilon'_{c_f} = G_{\theta_{\text{dm}} + \Delta\theta_{\text{dm}}}(x_t, c_f, t)$ ;
7    $\epsilon'_{c_\emptyset} = G_{\theta_{\text{dm}} + \Delta\theta_{\text{dm}}}(x_t, c_\emptyset, t)$ ;
8    $\mathcal{L}_{\text{cor}}(c_f, \Delta\theta_{\text{dm}}) = \text{Avg}((\epsilon_{c_f} - \epsilon_{c_\emptyset}) \cdot (\epsilon'_{c_f} - \epsilon'_{c_\emptyset}))$ ;
9    $\eta_i = \frac{\|\mathcal{L}_{\text{cor}}(c_{1,f}, \Delta\theta_{\text{dm}})\|_2}{\|\mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{\text{dm}})\|_2}$ ;
10   $\mathcal{L}_{\text{mom}}^n = \alpha \mathcal{L}_{\text{mom}}^{n-1} + (1 - \alpha) \sum_{i=1}^N \eta_i \mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{\text{dm}})$ ;
11  if  $\mathcal{L}_{\text{mom}}^n \leq \tau$  then
12    | break;
13  end
14   $\min_{\Delta\theta_{\text{dm}}} \mathcal{L}_{\text{G-CiRs}} = \sum_{i=1}^N \eta_i \mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{\text{dm}}) + \lambda \|\Delta\theta_{\text{dm}}\|_p$ 
15 end
```

C. Style removal

Reference images x_0 . Several visual examples generated by the original diffusion model are provided in Fig. 8.

Evaluation for G-CiRs. Additional visual examples produced with sentences³ containing artist names as prompts for various unlearned DMs are shown in Figs. 9~11.

Other experimental results. 1) The more detailed results of Tab. 4 are shown in Tab. 5. 2) Furthermore, we perform an ablation study on the threshold τ , which controls the moment of early stopping. The detailed experimental results are shown in Tab. 6. Observations reveal that G-CiRs achieves optimal erasing performance for ‘Van Gogh,’ ‘Picasso,’ and ‘Cezanne’ at τ values of -5e-4, 0, and 0, respectively.

D. Object Removal

Evaluation Settings: We employ the pre-trained ResNet50 [14] as the object classification network. We analyze nine classes within Imagenette [17], excluding the ‘cassette player’ category due to ResNet50’s classification accuracy falling below 50% on this class. During evaluation, we generate 250 images per class, *i.e.*, 50 seeds for each concept, and 5 images for each seed.

D.1. Evaluation for G-CiRs

For the *single concept erasure*, one category is chosen as the forgotten concept, while others serve as evaluation concepts. The quantitative results are presented in Tab. 7. On one hand, our G-CiRs achieves optimal performance in terms of ACC metric across nine object categories. On the other hand, the proposed G-CiRs induces fewer modifications to model weights. Additionally, the qualitative results are presented in Figs. 12~14.

³https://github.com/rohitgandikota/erasing/blob/main/data/art_prompts.csv

Algorithm 2: SepME.

Input: The diffuser $G(\cdot)$, the frozen weights θ_{dm} , the weight increment $\Delta\theta_{\text{dm}}$, the N forgotten concepts $c_{i,f} \in c_f$, the blank prompt c_\emptyset , the inference dataset $x_0 \in D$, the noise schedule $\bar{\alpha}_t$, the hyperparameters λ and β .

Output: The fine-tuned weight increments $\Delta\theta_{i \in [1,N], \text{dm}}$.

```
1 for  $c_{i,f} \in c_f$  do
2    $A = [c_\emptyset^\top; c_{1,f}^\top; \dots; c_{i-1,f}^\top; c_{i+1,f}^\top; \dots; c_{N,f}^\top]^\top$ ;
3   Obtain solutions  $S_p, A \otimes S_p = 0$ ;
4   for  $\Delta\theta_{\text{to}_k} \in \Delta\theta_{i, \text{dm}}$  do
5     Initialize  $w \in \mathcal{W}$  to a zero matrix;
6      $\Delta\theta_{\text{to}_k} = (w \otimes (\beta S_p^\top))^\top$ 
7   end
8   for  $\Delta\theta_{\text{to}_v} \in \Delta\theta_{i, \text{dm}}$  do
9     Initialize  $w \in \mathcal{W}$  to a zero matrix;
10     $\Delta\theta_{\text{to}_v} = (w \otimes (\beta S_p^\top))^\top$ 
11  end
12 end
13 for  $n, x_0 \in D$  do
14   Randomly select a sampling step  $t$ ;
15    $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \epsilon \in \mathcal{N}(0, \mathbf{I})$ ;
16    $\epsilon_{c_f} = G_{\theta_{\text{dm}}}(x_t, c_f, t)$ ;
17    $\epsilon_{c_\emptyset} = G_{\theta_{\text{dm}}}(x_t, c_\emptyset, t)$ ;
18   for  $c_{i,f} \in c_f$  do
19      $\epsilon'_{c_{i,f}} = G_{\theta_{i, \text{dm}} + \Delta\theta_{\text{dm}}}(x_t, c_{i,f}, t)$ ;
20      $\epsilon'_{c_\emptyset} = G_{\theta_{i, \text{dm}} + \Delta\theta_{\text{dm}}}(x_t, c_\emptyset, t)$ ;
21      $\mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{\text{dm}}) = \text{Avg}((\epsilon_{c_{i,f}} - \epsilon_{c_\emptyset}) \cdot (\epsilon'_{c_{i,f}} - \epsilon'_{c_\emptyset}))$ ;
22      $\eta_i = \frac{\|\mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{\text{dm}})\|_2}{\|\mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{\text{dm}})\|_2}$ ;
23   end
24    $\mathcal{L}_{\text{mom}}^n = \alpha \mathcal{L}_{\text{mom}}^{n-1} + (1 - \alpha) \sum_{i=1}^N \eta_i \mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{\text{dm}})$ ;
25   if  $\mathcal{L}_{\text{mom}}^n \leq \tau$  then
26     break;
27   end
28    $\min_{\mathcal{W}} \mathcal{L}_{\text{SepME}} = \sum_{i=1}^N \eta_i \mathcal{L}_{\text{cor}}(c_{i,f}, \Delta\theta_{i, \text{dm}}) + \lambda \|\mathcal{W}\|_p$ 
29 end
```

D.2. Evaluation for SepME

Separate optimization. Next, we individually train $\Delta\theta_{1, \text{dm}}$, $\Delta\theta_{2, \text{dm}}$, and $\Delta\theta_{3, \text{dm}}$ before combining them to simultaneously erase multiple concepts. The experimental results are detailed in Tab. 8. It is apparent that all unlearning methods effectively erase forgotten concepts. However, AbConcept exhibits shortcomings in restoring a forgotten concept. For example, Abconcept under the setting $\theta_{\text{dm}} + \sum_{i=2}^3 \Delta\theta_{i, \text{dm}}$ only perform 74% classification accuracy for ‘chain saw’, respectively. In contrast, our SepME₁ in Tab. 8 demonstrates nearly perfect recovery of erased concepts. This emphasizes the effectiveness of our SepME in restoring forgotten concepts and the feasibility of separately optimizing $\Delta\theta_{i, \text{dm}}$.

Iterative-concept erasure. The iterative concept erasure implies that each erasure step t can only utilize knowledge of the previously forgotten concepts $c_{<t, f}$. To realize this erasure, we avoid sequentially fine-tuning model weights, as restoring the early weights $\Delta\theta_{i, \text{dm}}$ inevitably affects the erasure performance of $\Delta\theta_{>i, \text{dm}}$. Inspired by the success of separately optimizing $\Delta\theta_{1 \sim 3, \text{dm}}$, we endeavor to implement iterative concept erasure through this setting. In the initial unlearning step, we employ G-CiRs to fine-tune all parameters of cross-attention modules in DMs. This enables better unlearning of the forgotten concept with fewer weight modifications. In each subsequent unlearning step t , we recalculate S_p using $c_{<t, f}$ to construct the weight increments of to_k and to_v layers and fine-tune these increments to erase the concept $c_{t, f}$. The experimental results presented as SepME₂ in Tab. 8 demonstrate comparable performance to SepME₁. Overall, SepME can effectively



Figure 8. Visual examples produced by original diffusion models.

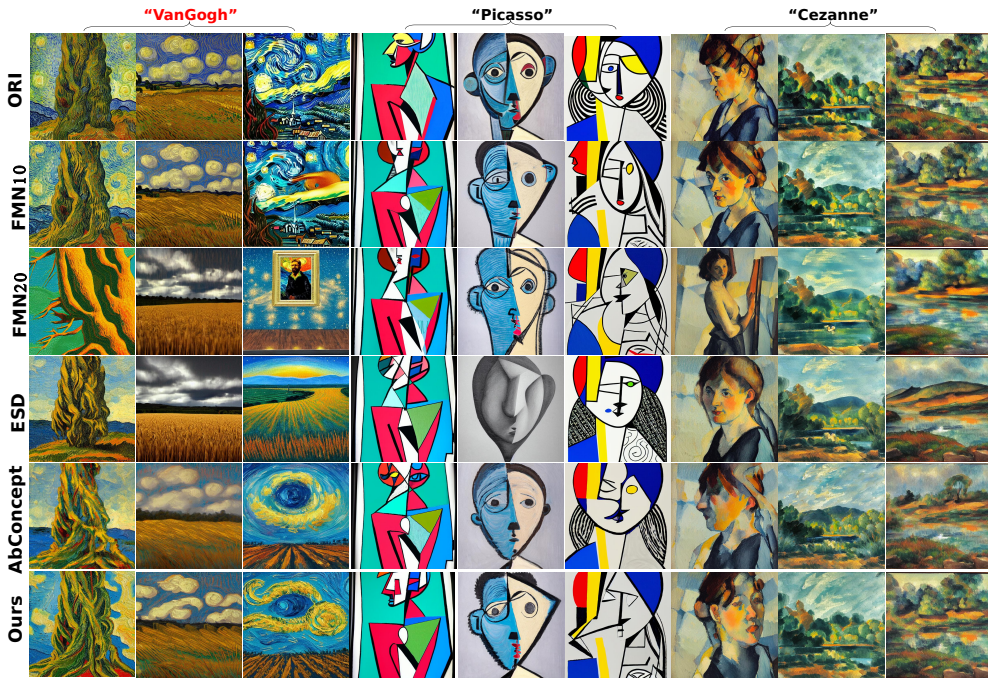


Figure 9. Qualitative comparison among various unlearning techniques for DMs with ‘Van Gogh’ as the erased concept.

achieve iterative concept erasure and concept restoration.

E. Other details.

We omit the consideration of iterative erasure, where multiple concepts are erased at each step, as the weight increments for erasing various concept can be optimized separately.

Cosine Function: We explored the use of the cosine function as an alternative to \mathcal{L}_{cor} and assessed its performance across various hyperparameters and learning rates. However, this approach did not yield satisfactory results. We attribute this to the significant prediction gap between model samples, *i.e.*, the normalization of sample constraints affects the optimization direction.

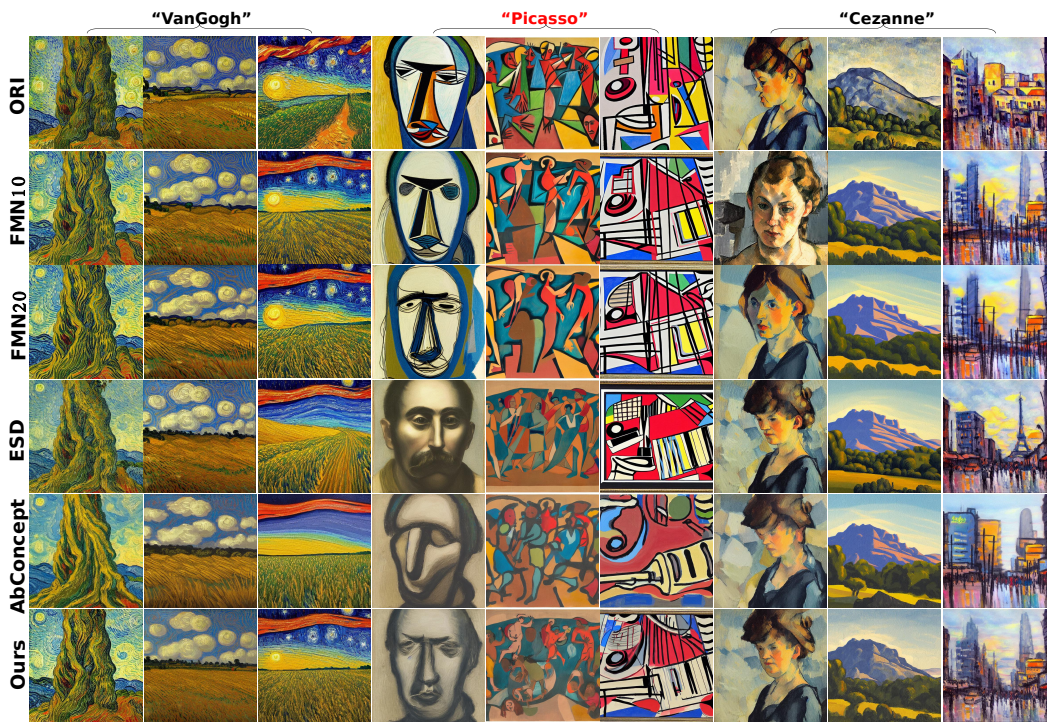


Figure 10. Qualitative comparison among various unlearning techniques for DMs with 'Picasso' as the erased concept.

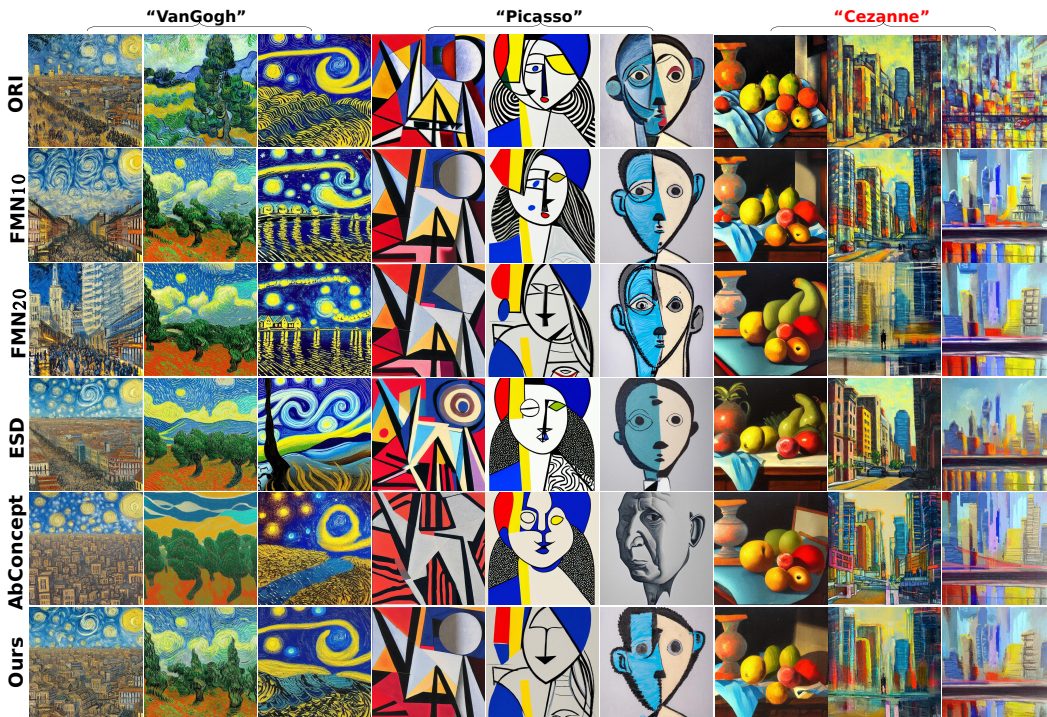


Figure 11. Qualitative comparison among various unlearning techniques for DMs with 'Cezanne' as the erased concept.

Table 5. Quantitative results of SepME when separately fine-tuning $\Delta\theta_{1\sim 3, \text{dm}}$. The concepts $c_{1,f}$, $c_{2,f}$, and $c_{3,f}$ correspond to ‘Van Gogh’, ‘Cezanne’, and ‘Picasso’, respectively. In SepME₁, $\Delta\theta_{1\sim 3, \text{dm}}$ are separately optimized when all forgotten concepts are known. SepME₂ follows the mode of iterative concept erasure. In other words, the t -th erasure step only possesses knowledge (Kn) of the previously forgotten concepts. ‘xattn’ means that we employ all layers of cross attention modules.

θ_{dm}	SepME ₁ -(ACC/LPIPS)- $Kn = [c_{1,f}; c_{2,f}; c_{3,f}]$						
	$+\Delta\theta_{1, \text{dm}}$	$+\Delta\theta_{2, \text{dm}}$	$+\Delta\theta_{3, \text{dm}}$	$+\sum_{i=1}^2 \Delta\theta_{i, \text{dm}}$	$+\sum_{i \in \{1,3\}} \Delta\theta_{i, \text{dm}}$	$+\sum_{i=2}^3 \Delta\theta_{i, \text{dm}}$	$+\sum_{i=1}^3 \Delta\theta_{i, \text{dm}}$
VG	0.356/0.364	1.000/0.182	0.908/0.182	0.308/0.365	0.308/0.365	<u>0.956/0.181</u>	0.308/0.364
CE	1.000/0.144	0.320/0.304	1.000/0.144	0.320/0.304	<u>1.000/0.144</u>	0.288/0.304	0.320/0.304
PC	1.000/0.185	1.000/0.185	0.000/0.460	<u>1.000/0.185</u>	0.000/0.460	0.000/0.460	0.000/0.460
Others	0.957/0.194	0.891/0.243	0.955/0.184	0.855/0.253	0.957/0.211	0.845/0.248	<u>0.803/0.260</u>
Kn	SepME ₂ -(ACC/LPIPS)						
	$[c_{1,f}]$	$[c_{1,f}; c_{2,f}]$	$[c_{1,f}; c_{2,f}; c_{3,f}]$	-	-	-	-
VG	0.228/0.363	0.956/0.182	0.956/0.182	0.228/0.363	0.228/0.363	<u>0.956/0.182</u>	0.228/0.363
CE	1.000/0.182	0.000/0.404	1.000/0.172	0.000/0.376	<u>1.000/0.177</u>	0.000/0.411	0.052/0.394
PC	0.964/0.150	1.000/0.144	0.288/0.277	<u>0.964/0.150</u>	0.108/0.292	0.288/0.277	0.108/0.292
Others	0.994/0.173	0.925/0.196	0.986/0.185	0.946/0.211	0.970/0.194	0.924/0.231	<u>0.916/0.235</u>
θ_{dm}	Abconcept-(ACC/LPIPS)-xattn						
VG	0.000/0.469	0.728/0.257	0.684/0.227	0.000/0.472	0.000/0.469	<u>0.636/0.277</u>	0.000/0.487
CE	0.820/0.217	0.036/0.360	0.752/0.194	0.288/0.299	<u>0.680/0.234</u>	0.252/0.288	0.144/0.349
PC	0.900/0.190	0.852/0.251	0.000/0.397	<u>0.600/0.277</u>	0.152/0.304	0.200/0.334	0.052/0.364
Others	0.952/0.185	0.923/0.215	0.971/0.173	0.861/0.227	0.891/0.194	0.903/0.206	<u>0.783/0.247</u>
θ_{dm}	Abconcept-(ACC/LPIPS)-(to_k;to_v)						
VG	0.592/0.329	0.820/0.215	0.728/0.271	0.320/0.383	0.184/0.378	<u>0.272/0.398</u>	0.184/0.411
CE	1.000/0.193	0.200/0.266	0.800/0.222	0.252/0.274	<u>0.600/0.230</u>	0.052/0.307	0.000/0.311
PC	1.000/0.142	0.964/0.141	0.216/0.306	<u>0.856/0.147</u>	0.320/0.316	0.072/0.347	0.144/0.356
Others	0.973/0.178	0.917/0.208	0.952/0.191	0.889/0.223	0.937/0.197	0.860/0.221	<u>0.805/0.248</u>
θ_{dm}	G-CiRs-(ACC/LPIPS)-xattn						
VG	0.228/0.363	0.908/0.180	0.908/0.276	0.184/0.386	0.092/0.382	<u>0.544/0.276</u>	0.184/0.410
CE	1.000/0.154	0.000/0.363	1.000/0.189	0.108/0.261	<u>0.716/0.178</u>	0.216/0.313	0.036/0.331
PC	1.000/0.175	1.000/0.204	0.052/0.358	<u>1.000/0.220</u>	0.452/0.247	0.200/0.280	0.200/0.303
Others	0.994/0.173	0.912/0.222	0.950/0.179	0.940/0.202	0.969/0.195	0.890/0.224	<u>0.843/0.268</u>
θ_{dm}	G-CiRs-(ACC/LPIPS)-(to_k;to_v)						
VG	0.344/0.356	0.908/0.209	0.820/0.236	0.228/0.378	0.228/0.381	<u>0.320/0.344</u>	0.136/0.397
CE	0.892/0.155	0.356/0.263	0.752/0.219	0.152/0.282	<u>0.900/0.210</u>	0.152/0.295	0.152/0.334
PC	1.000/0.184	1.000/0.192	0.148/0.290	<u>0.832/0.164</u>	0.252/0.293	0.144/0.319	0.072/0.346
Others	0.944/0.193	0.908/0.212	0.920/0.203	0.873/0.237	0.938/0.220	0.848/0.258	<u>0.791/0.272</u>

Table 6. Ablation study to investigate the influence of the hyperparameter τ on unlearning performance.

Erased (G-CiRs)	τ	Van Gogh	Picasso	Cezanne	Others	$\ \Delta\theta_{\text{dm}}\ _p \downarrow$
Van Gogh	1e-3	0.684/0.310	1.000/0.180	1.000/0.140	1.000/0.159	24.28
	5e-4	0.592/0.334	1.000/0.178	1.000/0.140	1.000/0.159	27.43
	0.	0.456/0.353	1.000/0.171	1.000/0.145	1.000/0.166	36.24
	-5e-4	0.228/0.363	1.000/0.175	1.000/0.154	0.992/0.173	43.17
	-1e-3	0.092/0.419	1.000/0.195	0.892/0.153	0.976/0.176	51.63
Picasso	5e-4	0.956/0.173	1.000/0.181	1.000/0.144	1.000/0.154	1.522
	5e-5	0.924/0.239	0.084/0.319	1.000/0.173	0.976/0.161	27.58
	0.	0.908/0.276	0.052/0.358	1.000/0.189	0.952/0.179	36.90
	-5e-4	0.044/0.573	0.000/0.567	0.000/0.426	0.812/0.325	65.03
	-1e-3	0.000/0.604	0.000/0.578	0.036/0.516	0.696/0.379	67.89
Cezanne	1e-4	0.772/0.216	1.000/0.182	0.276/0.071	0.976/0.179	23.51
	5e-5	0.772/0.242	1.000/0.215	0.072/0.355	0.960/0.195	30.67
	0.	0.908/0.180	1.000/0.204	0.000/0.363	0.912/0.222	36.55
	-5e-5	0.636/0.398	0.852/0.228	0.000/0.495	0.904/0.247	41.20
	-1e-4	0.272/0.458	0.852/0.251	0.000/0.543	0.832/0.281	45.75

Table 7. Quantitative results of the single concept erasure. i in FMN_i represents the iteration step.

ACC/LPIPS	ORI	FMN ₂₀		FMN ₅₀		AbConcept		G-CiRs	
		Erased	Others	Erased	Others	Erased	Others	Erased	Others
chain saw	0.96/0.	0.84/0.241	0.903/0.168	0.00/0.420	0.773/0.269	0.28/0.325	0.833/0.188	0.18/0.311	0.825/0.214
church	0.84/0.	0.84/0.243	0.893/0.167	0.06/0.440	0.870/0.187	0.30/0.345	0.870/0.183	0.16/0.255	0.863/0.184
gas pump	0.80/0.	0.20/0.256	0.930/0.170	0.02/0.379	0.835/0.263	0.20/0.359	0.905/0.186	0.12/0.341	0.893/0.178
tench	0.88/0.	0.26/0.412	0.873/0.178	0.00/0.462	0.818/0.256	0.08/0.403	0.870/0.173	0.08/0.381	0.878/0.177
garbage truck	0.94/0.	0.48/0.229	0.875/0.175	0.04/0.481	0.735/0.314	0.58/0.293	0.820/0.199	0.40/0.339	0.815/0.213
english springer	1.00/0.	0.82/0.221	0.873/0.167	0.02/0.364	0.810/0.256	0.14/0.260	0.873/0.232	0.00/0.292	0.833/0.232
golf ball	1.00/0.	0.86/0.277	0.850/0.171	0.44/0.420	0.765/0.222	0.00/0.459	0.845/0.231	0.32/0.345	0.875/0.218
parachute	0.98/0.	0.54/0.429	0.858/0.172	0.02/0.493	0.770/0.258	0.40/0.436	0.853/0.202	0.28/0.448	0.823/0.204
french horn	1.00/0.	0.22/0.377	0.853/0.185	0.04/0.429	0.665/0.291	0.00/0.435	0.823/0.185	0.00/0.452	0.818/0.188
average	0.93/0.	0.63/0.298	0.879/0.173	0.07/0.432	0.782/0.257	0.22/0.368	0.855/0.199	0.17/0.352	0.847/0.201
$\ \Delta\theta_{dm}\ _p \downarrow$	-	89.30		254.9		150.6		37.92	

Table 8. Quantitative results of SepME when separately fine-tuning $\Delta\theta_{1\sim 3, dm}$. $\Delta\theta_{1, dm}$, $\Delta\theta_{2, dm}$, $\Delta\theta_{3, dm}$ are optimizable weights for erasing ‘chain saw’, ‘gas pump’ and ‘garbage truck’, respectively. In SepME₁, $\Delta\theta_{1\sim 3, dm}$ are separately optimized when all forgotten concepts are known. SepME₂ follows the mode of iterative concept erasure. In other words, the t -th erasure step only possesses knowledge of the previously forgotten concepts. ‘xattn’ means that we employ all layers of cross attention modules.

ACC/LPIPS	SepME ₁ (to_k, to_v)								
	θ_{dm}	$+\Delta\theta_{1, dm}$	$+\Delta\theta_{2, dm}$	$+\Delta\theta_{3, dm}$	$+\sum_{i=1}^2 \Delta\theta_{i, dm}$	$+\sum_{i \in \{1, 3\}} \Delta\theta_{i, dm}$	$+\sum_{i=2}^3 \Delta\theta_{i, dm}$	$+\sum_{i=1}^3 \Delta\theta_{i, dm}$	
chain saw	0.080/0.311	0.940/0.177	0.940/0.182	0.140/0.314	0.120/0.312	0.960/0.184	0.120/0.314		
gas pump	0.680/0.142	0.300/0.308	0.680/0.142	0.260/0.307	0.680/0.142	0.360/0.298	0.340/0.298		
garbage truck	0.920/0.157	0.920/0.157	0.000/0.410	0.920/0.157	0.000/0.411	0.000/0.411	0.000/0.411		
Others	0.900/0.194	0.932/0.166	0.892/0.208	0.912/0.196	0.840/0.243	0.863/0.221	0.833/0.247		
		AbConcept (xattn)							
chain saw	0.280/0.325	0.940/0.192	0.920/0.203	0.340/0.328	0.220/0.344	0.740/0.222	0.080/0.359		
gas pump	0.560/0.175	0.200/0.359	0.520/0.182	0.520/0.310	0.460/0.212	0.560/0.299	0.400/0.331		
garbage truck	0.900/0.180	0.860/0.171	0.580/0.293	0.800/0.202	0.620/0.238	0.720/0.221	0.580/0.264		
Others	0.867/0.192	0.907/0.187	0.853/0.201	0.843/0.234	0.827/0.246	0.843/0.235	0.687/0.309		
		SepME ₂ (to_k, to_v)							
chain saw	0.220/0.311	0.940/0.176	0.940/0.181	0.220/0.311	0.120/0.321	0.980/0.181	0.140/0.333		
gas pump	0.540/0.187	0.380/0.214	0.680/0.151	0.260/0.230	0.580/0.187	0.480/0.200	0.280/0.224		
garbage truck	0.840/0.203	0.800/0.162	0.000/0.410	0.660/0.185	0.000/0.406	0.000/0.396	0.000/0.387		
Others	0.870/0.234	0.910/0.174	0.890/0.208	0.847/0.239	0.743/0.275	0.883/0.221	0.710/0.278		

Remove "Church"



Figure 12. Qualitative comparison among various unlearning techniques for DMs with 'church' as the erased concept.

Remove "English Springer"



Figure 13. Qualitative comparison among various unlearning techniques for DMs with 'English springer' as the erased concept.

Remove "Parachute"



Figure 14. Qualitative comparison among various unlearning techniques for DMs with 'parachute' as the erased concept.