
Choose Your Anchor Wisely: Effective Unlearning Diffusion Models via Concept Reconditioning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large-scale conditional diffusion models (DMs) have demonstrated exceptional
2 ability in generating high-quality images from textual descriptions, gaining
3 widespread use across various domains. However, these models also carry the risk
4 of producing harmful, sensitive, or copyrighted content, creating a pressing need to
5 remove such information from their generation capabilities. While retraining from
6 scratch is prohibitively expensive, machine unlearning provides a more efficient
7 solution by selectively removing undesirable knowledge while preserving utility. In
8 this paper, we introduce **C**oncept **R**Econditioning (**CORE**), a simple yet effective
9 approach for unlearning diffusion models. Similar to some existing approaches,
10 CORE guides the noise predictor conditioned on forget concepts towards an anchor
11 generated from alternative concepts. However, CORE introduces key differences
12 in the choice of anchor and retain loss, which contribute to its enhanced perform-
13 mance. We evaluate the unlearning effectiveness and retainability of CORE on
14 UnlearnCanvas. Extensive experiments demonstrate that CORE surpasses state-of-
15 the-art methods including its close variants and achieves near-perfect performance,
16 especially when we aim to forget multiple concepts. More ablation studies show
17 that CORE’s careful selection of the anchor and retain loss is critical to its superior
18 performance.

19 1 Introduction

20 In recent years, large-scale text-to-image generative models, especially Diffusion Models (DM), have
21 made remarkable advancements in artificial intelligence by exhibiting an unprecedented ability to
22 create high-resolution, high-quality images from text descriptions (Sohl-Dickstein et al., 2015; Ho
23 et al., 2020; Rombach et al., 2022). The versatility and accessibility of diffusion models have led
24 to their widespread adoption across various industries (Croitoru et al., 2023; Kazerouni et al., 2023;
25 Yang & Hong, 2022; Xu et al., 2022).

26 Despite their broad utility, diffusion models come with inherent risks due to their extensive training
27 on diverse datasets. These models have the potential to generate inappropriate, harmful, or legally
28 sensitive content. For example, Stable Diffusion can produce images that involve pornography,
29 malign stereotypes, and gender and race biases based on the embedded prejudices in their training
30 data, even conditional on non-harmful prompts (Birhane et al., 2021; Schramowski et al., 2023;
31 Larrazabal et al., 2020). They can memorize and reproduce realistic yet inappropriate depictions
32 of individuals without their consent, posing huge privacy risks (Somepalli et al., 2023a,b; Carlini
33 et al., 2023). They can also create misleading or harmful media involving real individuals, such as
34 deepfakes (Mirsky & Lee, 2021). Moreover, they can mimic potentially copyrighted content and
35 replicate styles of real artists, raising legal concerns related to copyright infringement and intellectual

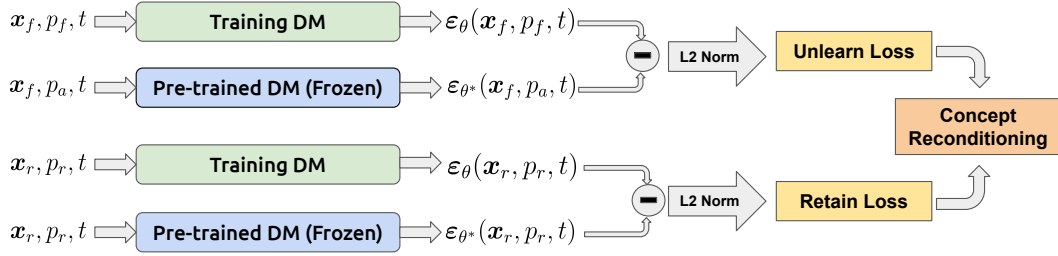


Figure 1: Overview of Concept Reconditioning. p_f, p_r, p_a are the concepts targeted to be forgotten (i.e., *forget concepts*), to be remembered (i.e., *retain concepts*), and to guide unlearning (i.e., *alternative concepts*), respectively. t is the number of steps in the denoising process and is uniformly sampled within $[0, T]$, where T denotes the total number of denoising steps in diffusion models. ε_θ is the noise predictor function we aim to optimize, while ε_{θ^*} is the noise predictor in the pre-trained diffusion models.

36 property rights, as well as undermining artistic originality (Shan et al., 2023; Roose, 2022; Liu, 2022;
 37 Popli, 2022; Scenario, 2022; Brittan, 2023).

38 To address these concerns, legislative frameworks such as the European Union’s General Data
 39 Protection Regulation (GDPR) (Mantelero, 2013; Voigt & Von dem Bussche, 2017) and the US’s
 40 California Consumer Privacy Act (CCPA) (CCPA, 2018) have established the Right to be Forgotten.
 41 These laws mandate that applications must support the deletion of personal information contained in
 42 training samples upon user request. Consequently, there is a pressing need for effective methods to
 43 mitigate these risks by enabling diffusion models to **unlearn** such undesirable content, ensuring that
 44 their deployment is both responsible and aligned with societal values.

45 A straightforward method is to retrain the model from scratch using a filtered dataset devoid of
 46 inappropriate content. However, this approach is computationally intensive and often impractical due
 47 to the enormous resources required. For instance, training Stable Diffusion 2.0 on a filtered image
 48 set (Schuhmann et al., 2022; Rombach & Esser, 2022) demands approximately 150,000 GPU hours
 49 on 256 A100 GPUs. Early attempts to unlearn large-scale generative models include decoding-time
 50 guidance and post-generation filtering (Rando et al., 2022; Schramowski et al., 2023); however, these
 51 methods do not modify the model weights and can be easily bypassed during deployment. Recent
 52 research has pivoted towards more robust fine-tuning-based unlearning approaches that modify a
 53 model’s weights to effectively forget specific undesirable elements (Gandikota et al., 2023; Fan et al.,
 54 2023; Heng & Soh, 2024; Kumari et al., 2023; Wu et al., 2024; Zhang et al., 2024a; Wu & Harandi,
 55 2024; Li et al., 2024b). These methods aim to steer the noise predictor in diffusion models away from
 56 the target concepts intended to be forgotten by efficiently fine-tuning a small fraction of parameters.

57 In this work, we propose **Concept REconditioning (CORE)**, a novel, simple, but effective unlearning
 58 method for diffusion model. This method leverages a fixed, non-trainable noise to guide the
 59 unlearning process, circumventing the need for dual noise predictors or the use of Gaussian noise as a
 60 target. CORE specifically alters the noise prediction mechanism for the target images conditioned on
 61 concepts in the forget set (i.e., *forget concepts*), aligning them closer to concepts in the retain set (i.e.,
 62 *retain concepts*), thereby blurring the distinction between correctly generated images from forget
 63 concepts and incorrectly generated ones from retain concepts. We position CORE within a more
 64 general framework of *Concept Erasing*, and compare our method with other baselines that fit into this
 65 framework. Despite its simplicity, we demonstrate its superiority over existing methodologies through
 66 rigorous testing on the UnlearnCanvas framework, and show CORE excels in overall performance
 67 including unlearning ability, retainability, and generalization ability, especially when we aim to forget
 68 multiple concepts.

69 Our contributions are summarized as follows.

- 70 • We introduce **Concept REconditioning (CORE)** as a new efficient and effective unlearning
 71 method on diffusion models, and position it in a broader conceptual framework of concept erasing.
- 72 • Extensive empirical validations on UnlearnCanvas showcase that CORE significantly outperforms
 73 existing baselines, achieving nearly perfect scores and setting new state-of-the-arts for the over-
 74 all performance in unlearning diffusion models on UnlearnCanvas. CORE also shows strong
 75 capabilities of generalization in unlearning styles.

76 • Ablation studies highlight the benefits of using a fixed, non-trainable target noise over other
 77 methods. Additionally, our findings emphasize the superiority of one-to-one concept reconditioning
 78 over other schemes of selecting reconditioning concepts.

79 2 Preliminaries

80 **Machine Unlearning.** Machine Unlearning (MU) refers to the process of systematically removing
 81 the influence of specific data points from a trained machine learning model, ensuring that the model
 82 forgets information as if the data points were never included in its training set. In this context, let
 83 \mathcal{D} represent the training dataset, and let $\mathcal{D}_f \subset \mathcal{D}$ denote the forget set, the subset of data that needs
 84 to be unlearned. The retain set, denoted as $\mathcal{D}_r \subset \mathcal{D}$, is the complement of the forget set. The goal
 85 of machine unlearning is to produce a new model that closely approximates the performance of
 86 retraining from scratch on \mathcal{D}_r while also ensuring that the model does not retain any knowledge
 87 of \mathcal{D}_f . Unlearning has traditionally been explored in the context of classification models, where
 88 the model aims to either forget the influence of specific classes of data or forget some random
 89 samples (Cao & Yang, 2015; Bourtoule et al., 2021). In recent developments, machine unlearning
 90 has been extended to large generative models, where the model must unlearn specific objectives to
 91 ensure that certain generated outputs, such as sensitive, private, copyrighted, or harmful content, will
 92 not be generated.

93 **Unlearning Diffusion Models.** Diffusion models are a class of generative models that have gained
 94 significant attention for their ability to generate high-quality images. They work by transforming
 95 data distributions through T forward and reverse steps, gradually adding noise to the data and then
 96 learning to reverse this process to generate new samples. Mathematically, this can be described by a
 97 series of noisy images $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T \in \mathbb{R}^d$, where \mathbf{x}_0 is the original image, and \mathbf{x}_T is the Gaussian
 98 noise. Latent Diffusion Model (LDM) (Rombach et al., 2022) first compresses high-dimensional
 99 pixel-based data into a low-dimensional latent space using an encoder \mathcal{E} . It then simulates the
 100 diffusion process on the space of latent variables $\mathbf{z} = \mathcal{E}(\mathbf{x})$ and reconstructs the image through a
 101 decoder \mathcal{D} . For notational simplicity, we do not differentiate between latent variables and pixel-based
 102 data, denoting both as \mathbf{x} . In this context, let $\varepsilon_\theta(\mathbf{x}_t, p)$ represent the noise estimator parameterized by
 103 θ , where \mathbf{x}_t is the noisy observation at step t , and p is a conditioning variable such as a class label or
 104 text description. The training objective of latent diffusion models is the mean squared error (MSE)
 105 between the predicted noise and the true noise across all diffusion steps, expressed as:

$$\mathcal{L}_{\text{MSE}}(\theta) = \mathbb{E}_{p,t,\varepsilon \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \left[\|\varepsilon - \varepsilon_\theta(\mathbf{x}_t, p)\|_2^2 \right], \quad (1)$$

106 where p is sampled from a distribution over all prompts and t is sampled uniformly from $[0, T]$. Given
 107 a pre-trained latent diffusion model, the objective of unlearning this diffusion model is to ensure that
 108 harmful or sensitive content, such as depictions of nudity or violence, can no longer be produced by
 109 the model when prompted with the corresponding text descriptions. The challenge lies in balancing
 110 the removal of unwanted generations while preserving the model’s ability to generate high-quality,
 111 appropriate content for normal prompts. The most common unlearning process in diffusion models
 112 involves updating the noise estimator to ensure that harmful concepts associated with \mathcal{D}_f are no
 113 longer learned or reinforced during the reverse diffusion process. This form of unlearning, often
 114 referred to as “concept erasure”, is critical for ensuring the safe deployment of generative models in
 115 real-world applications. More details are included in Section 3.2.

116 3 Concept Reconditioning

117 In this section, we propose **C**oncept **R**Econditionng (**CORE**), a simple yet effective algorithm
 118 for unlearning in diffusion models. Our approach focuses on reconditioning the model’s learned
 119 representations by substituting concepts from the forget set with selected alternative concepts from
 120 the retain set. First, we introduce the objective function and key designs within. Then, we position it
 121 within the broader framework of *Concept Erasing* and compare it with similar algorithms in prior
 122 works to showcase its advantage.

123 3.1 Proposed Method

124 **Unlearn objective.** In the context of unlearning in diffusion models, we denote the noise predictor in
 125 Latent Diffusion Models by $\varepsilon_\theta(\mathbf{x}_t, p)$, where \mathbf{x}_t is the noisy version of the input image \mathbf{x}_0 at time

126 step t generated during the forward diffusion process, p is the prompt associated with the image (e.g.,
 127 “A cat in the style of Van Gogh”), and θ represents the model parameters. We use $\varepsilon_{\theta^*}(\mathbf{x}_t, p)$ and θ^*
 128 to denote the pre-trained diffusion model and its parameter. In CORE, we aim to recondition images
 129 from the forget set onto alternative concepts. This is achieved by aligning the noise estimator for
 130 images in the forget set, conditioned on their original concepts $p_f \in \mathcal{D}_f$, toward the ground truth
 131 noise estimator for the same image but conditioned on an alternative concept p_a . Mathematically, the
 132 unlearn objective function is formulated as

$$\mathcal{L}_f(\theta) := \mathbb{E}_{(p_f, \mathbf{x}_0) \sim \mathcal{D}_f, p_a \neq p_f, t} [\|\varepsilon_{\theta}(\mathbf{x}_t, p_f) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_a)\|_2^2], \quad (2)$$

133 where the expectation is taken over the concept-image pairs (p_f, \mathbf{x}_0) from the forget set, alternative
 134 concepts p_a different from p_f , and time steps t uniformly sampled from $[0, T]$. Intuitively, this
 135 process effectively weakens the association between the images and their original concepts in the
 136 model, steering it away from the initial pre-trained associations.

137 **Alternative concepts.** A key design choice in CORE is the selection of alternative concepts p_a
 138 in equation (2). In the unlearning objective, p_a acts as an anchor concept to recondition images
 139 from the forget set onto. Previous works typically use an empty string or a single base concept for
 140 p_a consistently across all concepts to be unlearned (Zhang et al., 2024c; Gandikota et al., 2023).
 141 In contrast, CORE adopts a different approach by pairing each forget concept p_f with a specific
 142 alternative concept p_a . Our pairing scheme imposes minimal restrictions: the alternative concept p_a
 143 does not necessarily have to come from the retain set; it can even be another forget concept different
 144 from p_f . In our implementation, when the number of concepts to forget is smaller than the number of
 145 retain concepts, we map each forget concept to a unique concept in the retain set, rather than using a
 146 single base concept for all forget concepts. Meanwhile, when the retain concepts are limited and there
 147 are more concepts to forget, we create a one-to-one mapping among the forget concepts themselves.
 148 This means that each forget concept p_f is paired with another forget concept p_a (where $p_a \neq p_f$) to
 149 serve as its alternative concept during unlearning. Empirically, we show that this one-to-one mapping
 150 strategy significantly outperforms methods that consistently use a base concept or randomly sample
 151 alternative concepts at each step.

152 **Retain objective and the full loss function.** To ensure the model continues generating high-quality
 153 images for the retain concepts, we introduce a retain loss to regularize the unlearning process.
 154 Traditionally, the retain loss is defined as the Mean Squared Error (MSE) between the noise prediction
 155 for the retain set and the Gaussian noise vector used to generate the noisy images, similar to the
 156 objective used in training a diffusion model (see equation 1). However, in CORE, rather than fine-
 157 tuning the noise predictions to match a Gaussian random vector, we instead align them with those
 158 generated by the pre-trained diffusion model itself. Mathematically, the retain objective is defined as

$$\mathcal{L}_r(\theta) := \mathbb{E}_{(p_r, \mathbf{x}_0) \sim \mathcal{D}_r, t} [\|\varepsilon_{\theta}(\mathbf{x}_t, p_r) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_r)\|_2^2], \quad (3)$$

159 where t is uniformly sampled in $[0, T]$ and (p_r, \mathbf{x}_0) are concept-image pairs sampled from the
 160 retain set. Using $\varepsilon_{\theta^*}(\mathbf{x}_t, p_r)$ as the target helps ensure the model does not deviate too far from its
 161 original capabilities, as it leverages the pre-trained model’s learned knowledge. Empirical results (see
 162 Section 4) demonstrate that aligning the noise predictions with $\varepsilon_{\theta^*}(\mathbf{x}_t, p_r)$, rather than the Gaussian
 163 noise, yields better performance. This improvement arises potentially because using the estimated
 164 noise from the pre-trained model reduces variance in the unlearned model and stabilize the training
 165 process. Interestingly, this phenomenon, where using estimated signals outperforms true signals, has
 166 also been observed in other domains in statistics (Robins et al., 1992; Henmi & Eguchi, 2004; Hitomi
 167 et al., 2008; Su et al., 2023).

168 Finally, the complete loss function in CORE combines both the unlearn and retain objectives:

$$\mathcal{L}(\theta) := \mathcal{L}_f(\theta) + \alpha \cdot \mathcal{L}_r(\theta), \quad (4)$$

169 where $\alpha > 0$ controls the regularization strength. Intuitively, CORE ensures that the model is
 170 steered away from generating images associated with forget concepts while preserving its overall
 171 performance on other concepts.

172 3.2 Rethinking Concept Erasing and Reconditioning

173 At first glance, our proposed objective might seem similar to existing methods for unlearning in
 174 diffusion models, as it also involves steering the error predictor on the forget set while keeping it

175 unchanged on the retain set. However, under closer scrutiny, Concept Reconditioning introduces
 176 several key distinctions that set it apart and enable it to outperform previous approaches. Take a
 177 broader view of the framework of unlearning diffusion models: unlearning methods for diffusion
 178 models that are based on fine-tuning the error predictor $\varepsilon_\theta(\mathbf{x}, p)$ can generally be categorized into
 179 two classes: ❶ Concept Erasing (CE): This method works by shifting the noise prediction network
 180 for images corresponding to the forget concepts towards an alternative noise distribution. Intuitively,
 181 by doing so, it directly acts on $\varepsilon_\theta(\mathbf{x}_t^f, p_f)$, where \mathbf{x}_t^f is the noisy observation for images in the forget
 182 set, and misleads them away. ❷ Image Relabeling (IR): In this approach, alternative images that do
 183 not match the forget concepts are selected, and the model is fine-tuned on the forget concepts paired
 184 with these mismatched images. The model directly acts on $\varepsilon_\theta(\mathbf{x}_t^r, p_f)$ where \mathbf{x}_t^r is the noisy images
 185 constructed from the retain set, and effectively overwrites the old knowledge with new associations,
 186 forcing it to forget by learning new, incorrect pairings. Mathematically, these two classes can be
 187 formulated as

$$\mathcal{L}_{\text{CE}}(\theta) := \lambda \cdot \mathbb{E}_{(p_f, \mathbf{x}_0) \sim \mathcal{D}_{f,t}} [\|\varepsilon_\theta(\mathbf{x}_t, p_f) - \mathbf{y}_{\text{CE}}\|_2^2], \quad (5)$$

$$\mathcal{L}_{\text{IR}}(\theta) := \lambda \cdot \mathbb{E}_{p_f \sim \mathcal{D}_f, \mathbf{x}_0 \sim \mathcal{D}_r, t} [\|\varepsilon_\theta(\mathbf{x}_t, p_f) - \mathbf{y}_{\text{IR}}\|_2^2]. \quad (6)$$

188 Here, $\lambda \in \{\pm 1\}$ controls the direction of the objective function. In the CE method, images are drawn
 189 from the forget set, while in IR, images come from the retain set. The **target noises** \mathbf{y}_{CE} and \mathbf{y}_{IR} can
 190 be either random vectors (e.g., Gaussian or Uniform) or derived from a trainable noise predictor.

191 Many existing unlearning methods fit within this framework. For example, Heng & Soh (2024)
 192 suggests $\lambda = -1$ and $\mathbf{y}_{\text{CE}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ in equation (5) in the unlearning objective, while proposing
 193 a surrogate objective with $\lambda = 1$ and $\mathbf{y}_{\text{IR}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ in equation (6). The former corresponds to
 194 a gradient ascent loss applied to the pre-training objective on forget concepts, while the surrogate
 195 objective simply mirrors the standard training loss applied to the forget concepts with retain images.
 196 Fan et al. (2023) takes \mathbf{y}_{CE} in equation (5) as a trainable noise predictor $\varepsilon_\theta(\mathbf{x}_t, p_a)$ where $p_a \neq p_f$ is
 197 an alternative concept coming from the retain set. Wu et al. (2024) also proposes this target noise, as
 198 well as suggesting an alternative with \mathbf{y}_{CE} as a uniformly distributed random vector. Kumari et al.
 199 (2023) takes \mathbf{y}_{IR} in equation (6) to be either a standard Gaussian random vector or the error predictor
 200 at the last iterate, evaluated at retain images paired with corresponding retain concepts. Even when
 201 the objective function appears divergent from this framework, as seen in Gandikota et al. (2023), it
 202 can still be decomposed into a linear combination of objective functions in the framework above (see
 203 Appendix C).

204 Although these prior works often include additional techniques such as weight decay (Heng & Soh,
 205 2024), saliency map (Fan et al., 2023), or even applying a monotonic function to the squared loss (Park
 206 et al., 2024), the backbone of their unlearning objectives can be positioned into this simple framework
 207 or its simple variants. Our method distinguishes itself from prior approaches by its simplicity. Unlike
 208 previous methods, CORE requires no auxiliary techniques, and simply optimizing the objective $\mathcal{L}(\theta)$
 209 in equation (4) achieves state-of-the-art results.

210 Another key distinction is that CORE uses a fixed, non-trainable noise predictor from the pre-trained
 211 diffusion model as the target noise. This fixed anchor provides a clearer target noise compared
 212 to a trainable network or a random vector with a fixed distribution (e.g., a uniformly distributed
 213 random vector). Let us compare the three types of target noises. With a random vector from a
 214 fixed distribution (Kumari et al., 2023; Heng & Soh, 2024), there is no guarantee that this manually
 215 designed random vector will effectively disrupt the noise predictor conditioned on the forget concepts.
 216 A trainable, non-fixed noise (Fan et al., 2023; Kumari et al., 2023; Wu et al., 2024) is unstable during
 217 the unlearning process, particularly when aiming to forget many concepts over a long training period,
 218 since this target may drift towards an undesired direction. While methods using trainable target noises
 219 include a retain term in their loss function, this retain objective directly influences $\varepsilon_\theta(\mathbf{x}_t^r, p_r)$ but not
 220 $\varepsilon_\theta(\mathbf{x}_t^f, p_r)$, where \mathbf{x}_t^r and \mathbf{x}_t^f are noisy observations from the retain and forget sets, respectively. In
 221 contrast, CORE’s use of a non-trainable target noise ensures that the noise predictor always learns
 222 from a reference “incorrect” noise estimator derived from the pre-trained model.

223 4 Experiments

224 In this section, we show CORE outperforms baselines on UnlearnCanvas (Zhang et al., 2024c).

225 4.1 Experiment Setup

226 **Dataset and Tasks.** UnlearnCanvas
 227 is a high-resolution stylized image
 228 dataset designed to evaluate diffusion
 229 model unlearning methods (Zhang
 230 et al., 2024c). The dataset consists
 231 of images across 50 unique styles and
 232 20 distinct objects, with 20 images for
 233 each style-object combination. Each
 234 image is labeled with both a style and
 235 an object, making it particularly well-
 236 suited for measuring the unlearning ef-
 237 fectiveness and the retainability both
 238 within a single domain and across do-
 239 mains. In this paper, we mainly fo-
 240 cus on style unlearning within the Un-
 241 learnCanvas dataset. We define three
 242 unlearning tasks, each progressively forgetting more styles: Forget01 (forgetting 1 style), Forget06
 243 (forgetting 6 styles), and Forget25 (forgetting 25 styles).

244 Models and Baselines.

245 We use a Stable Diffusion v1.5
 246 model (Rombach et al., 2022) to per-
 247 form the fine-tuning and unlearning,
 248 and we also use a vision Transformer
 249 (ViT-Large) (Dosovitskiy, 2020) on
 250 UnlearnCanvas for style and object
 251 classification. Before unlearning the
 252 model, the base Stable Diffusion
 253 model is fine-tuned on all images
 254 from UnlearnCanvas. After complet-
 255 ing the unlearning phase, we prompt
 256 the unlearned model to generate im-
 257 ages conditioned on concepts from
 258 both forget and retain sets. The vi-
 259 sion Transformer is then used to clas-
 260 sify the generated images and calcu-
 261 late the relevant metrics. We compare
 262 CORE with several state-of-the-art un-
 263 learning methods for diffusion mod-
 264 els, including ESD (Gandikota et al.,
 265 2023), SalUn (Fan et al., 2023), Ed-
 266 iff (Wu et al., 2024), CA-model and
 267 CA-noise (Kumari et al., 2023). See
 268 Appendix C for more details.

269 **Metrics.** Following Zhang et al.
 270 (2024c), we use Unlearning Accuracy
 271 (UA) to assess the unlearning effec-
 272 tiveness. UA is the percentage of
 273 images generated by the unlearned
 274 model, conditioned on the forget con-
 275 cepts, which are incorrectly classified
 276 by the vision Transformer. A higher
 277 UA indicates stronger unlearning ca-
 278 pabilities. We measure retainability



Figure 2: Generated images from the unlearned model. The first column is generated by the fine-tuned Stable Diffusion model before any unlearning. Other columns are generated by the model unlearned by our proposed method and five baseline methods. More images are included in Appendix D.

Figure 2: Generated images from the unlearned model. The first column is generated by the fine-tuned Stable Diffusion model before any unlearning. Other columns are generated by the model unlearned by our proposed method and five baseline methods. More images are included in Appendix D.

Algorithm	UA (↑)	IRA (↑)	CRA (↑)	SFID (↑)	Total (↑)
Forget01					
Original	0.00	100.00	96.67	100.00	296.67
Ediff	93.33	84.00	98.33	100.00	375.66
CA-model	96.67	80.00	92.78	100.00	369.45
CA-noise	100.00	100.00	96.11	100.00	396.11
SalUn	53.33	98.67	92.78	95.74	340.52
ESD	100.00	66.00	96.11	96.95	359.06
CORE (ours)	93.33	98.00	96.11	100.00	387.44
Forget06					
Original	0.00	100.00	98.33	100.00	298.33
Ediff	45.00	80.00	99.17	100.00	324.17
CA-model	85.00	81.67	88.33	88.09	343.09
CA-noise	85.00	91.67	85.83	92.46	354.96
SalUn	90.00	83.33	98.33	88.52	360.18
ESD	100.00	75.00	100.00	93.47	368.47
CORE (ours)	90.00	100.00	97.50	99.56	387.06
Forget25					
Original	1.20	96.54	95.29	100.00	293.03
Ediff	54.00	78.46	95.10	84.48	312.04
CA-model	68.60	78.85	95.69	81.73	324.87
CA-noise	47.20	86.15	90.59	82.09	306.03
SalUn	51.60	77.31	87.65	82.34	298.90
ESD	90.40	46.54	99.02	88.12	324.08
CORE (ours)	91.60	95.38	97.65	100.00	384.63

Table 1: Performance of CORE and five baseline methods using Stable Diffusion v-1.5 on Forget01, Forget06, and Forget25 in UnlearnCanvas. Unlearning accuracy, In-domain and cross-domain retain accuracy, and scaled FID value serve as main metrics and are summarized in Section 4.1. For details about the scaled FID value, see Appendix B. The best total score is highlighted in **bold**.

¹There are 60 styles in UnlearnCanvas dataset, but in its latest codebase only 50 styles are used. See <https://github.com/OPTML-Group/UnlearnCanvas>.

279 Cross-domain Retain Accuracy (CRA). IRA refers to the classification accuracy of generated images
 280 prompted with retain concepts, within the same domain (e.g., when forgetting “Van Gogh’s style”,
 281 an in-domain prompt might be “A painting in crayon style”). CRA measures accuracy for retain
 282 prompts across domains (e.g., for the same task, a cross-domain prompt might be “A painting of
 283 a cat,” specifying the object). Additionally, we evaluate the quality of generated images using the
 284 scaled FID (SFID) score, which maps the original FID score (Heusel et al., 2017) onto a 0–100 scale,
 285 where higher SFID values indicate better generation quality. We also present the summation of all
 286 four scores on a scale of 0-100, as a comprehensive measurement of the unlearning capacity and
 287 retainability. For more experimental details, see Appendix B.

288 4.2 Results

289 **CORE achieves the best overall performance.** In Table 1, we present the unlearning effectiveness
 290 and retainability of CORE compared to five baseline methods across Forget01, Forget06 and Forget25
 291 tasks from UnlearnCanvas. The “Original” row refers to the performance of the pre-trained model
 292 without any unlearning. On Forget01, CORE ranks second overall based on the total score. However,
 293 in the more challenging tasks Forget06 and Forget25, CORE consistently achieves the highest total
 294 score among all methods, with an increasing performance gap over the baseline methods. Notably,
 295 CORE is the only method that maintains strong performance as the size of the forget set grows. In
 296 the most difficult task, where 25 out of 50 concepts are targeted for forgetting, CORE achieves the
 297 highest unlearning accuracy, in-domain retain accuracy, and scaled FID score, while securing the
 298 second-best cross-domain retain accuracy. Compared to its close variants, ESD, CORE achieves
 299 similar unlearning accuracy but significantly outperforms in retainability, particularly in cross-domain
 300 tasks, due to the adoption of an additional retain loss. Compared to baseline methods that use a
 301 trainable noise predictor, such as SalUn and CA-model, CORE excels in forgetting more concepts
 302 due to the stability of its non-trainable target, which proves more reliable over longer unlearning
 303 periods. Figure 2 shows some generated images using CORE and five baseline methods.

304 **CORE shows better generalization ability in unlearning styles.**

305 We further investigate CORE’s ability to general-
 306 ize in unlearning styles, aiming
 307 to verify that CORE can effectively
 308 unlearn specific target styles, instead
 309 of simply overfitting to the training
 310 objects. To assess this, we train the
 311 model on only 10 objects for each for-
 312 get concept and then evaluate the un-
 313 learning accuracy on 10 unseen ob-
 314 jects. This tests the model’s ability
 315 to generalize beyond the specific ob-
 316 jects used during training. As shown
 317 in Table 2, CORE outperforms all baseline methods in terms of generalization ability.

Algorithm	UA (↑)	IRA (↑)	CRA (↑)	SFID (↑)	Total (↑)
Ediff	36.67	81.67	92.50	100.00	310.84
CA-model	85.00	83.33	96.67	86.67	351.67
CA-noise	81.67	91.67	87.50	87.62	348.46
SalUn	95.00	65.00	90.83	86.37	337.20
ESD	100.00	46.67	99.17	86.11	331.95
CORE (ours)	83.33	100.00	96.67	99.67	379.67

Table 2: Generalization ability of CORE and baseline methods using Stable Diffusion v-1.5 on Forget06 of UnlearnCanvas. Unlearning accuracy, In-domain and cross-domain retain accuracy, and scaled FID value serve as main metrics and are summarized in Section 4.1. The best total score is highlighted in **bold**.

319 **The role of non-trainable target noise.** A key design choice in CORE is the use of non-trainable
 320 target noise from the pre-trained diffusion model in both the unlearn and retain objectives. This is
 321 contrary to other approaches that use trainable noise predictors as targets in the unlearn loss and
 322 Gaussian noise vectors as targets in the retain loss. To isolate the specific effects of the non-trainable
 323 target noise, excluding the influence of auxiliary techniques like saliency maps, we evaluate several
 324 variants of CORE: ❶ We replace $\varepsilon_{\theta^*}(\mathbf{x}_t, p_a)$ with $\varepsilon_{\theta}(\mathbf{x}_t, p_a)$ in equation (2), where \mathbf{x}_t are noisy
 325 images from the forget set and p_a is the alternative concept. This variant mirrors the backbone of
 326 the unlearn loss used in SalUn (Fan et al., 2023). ❷ We replace $\varepsilon_{\theta^*}(\mathbf{x}_t, p_a)$ with a Gaussian noise
 327 ε in equation (2) and apply a negative sign to the unlearn loss. This variant follows the gradient
 328 ascent-based method, similar to the unlearn loss in CA-noise (Kumari et al., 2023). ❸ We replace
 329 $\varepsilon_{\theta^*}(\mathbf{x}_t, p_r)$ with a Gaussian noise ε in equation (3), where \mathbf{x}_t is noisy observations of images from
 330 the retain set. This variant is aligned with the retain loss employed in many baseline methods (Heng
 331 & Soh, 2024; Kumari et al., 2023; Wu et al., 2024). The results are shown in Table 3.

332 **Anchor Selection: How do we approach it?** Another key distinction between CORE and other
 333 baseline methods lies in how anchors p_a are selected in the unlearning objective (as defined in
 334 equation 2). In CORE, each forget concept p_f is paired with a distinct alternative concept. This

Unlearn Loss	Retain Loss	UA (↑)	IRA (↑)	CRA (↑)	SFID (↑)	Total (↑)
CORE	CORE	95.00	100.00	97.08	100.00	392.08
$\mathbb{E}\ \varepsilon_{\theta}(\mathbf{x}_t^f, p_f) - \varepsilon_{\theta^*}(\mathbf{x}_t^f, p_a)\ _2^2$	CORE	43.33	98.33	95.00	100.00	336.66
$-\mathbb{E}\ \varepsilon_{\theta}(\mathbf{x}_t^f, p_f) - \varepsilon\ _2^2$	CORE	85.00	61.67	60.00	79.48	286.15
CORE	$\mathbb{E}\ \varepsilon_{\theta}(\mathbf{x}_t^r, p_r) - \varepsilon\ _2^2$	83.33	93.33	96.67	99.92	373.25

Table 3: Performance of CORE and its variants on the Forget06 task from UnlearnCanvas. In each variant, one component of the loss function remains unchanged, while the non-trainable target noise in the other component is replaced with alternative approaches. Metrics and are summarized in Section 4.1. The best total score is highlighted in **bold**. Here, \mathbf{x}_t^f and \mathbf{x}_t^r are the noisy observations for images in the forget set and retain set, respectively; p_f, p_a, p_r correspond to forget concepts, alternative concepts, and retain concepts, respectively. ε denotes the standard Gaussian random vector used to generate \mathbf{x}_t^f . Here, we pair each forget concept with one distinct retain concept in all experiments above.

Scheme for reconditioned concepts	UA (↑)	IRA (↑)	CRA (↑)	SFID (↑)	Total (↑)
Default (one-to-one)	91.60	95.38	97.65	100.00	384.63
One base concept (all-to-one)	82.40	60.00	98.33	93.06	333.79
Five base concepts (five-to-one)	93.40	84.04	98.43	96.52	372.39
Random concept (one-to-all)	56.60	95.77	96.96	100.00	349.33
Random from five concepts (one-to-five)	56.40	95.58	97.75	99.78	349.51

Table 4: Comparison of different alternative concept selection schemes. All experiments are done in the Forget25 task from UnlearnCanvas. In CORE (referred to as "Default"), each forget concept is paired one-to-one with a distinct alternative concept. One base concept: all forget concepts are reconditioned onto a single base concept. Five base concepts: forget concepts are grouped into sets of five, with each group reconditioned to one base concept. Random concept: a random alternative concept is selected for each forget concept at every gradient step. Random from five concepts: each forget concept is paired with five alternative concepts, with one randomly sampled at each step. The best total score is highlighted in **bold**. Significant underperforming results are highlighted in green.

335 contrasts with other methods that recondition all forget concepts to a single base concept or the empty
336 string. To demonstrate the effectiveness of CORE’s one-to-one pairing, we compare different selection
337 schemes: One approach involves pairing each forget concept with a set of alternative concepts (or
338 even the entire retain set) and randomly sampling one at each gradient step to recondition the target
339 images. Another approach reconditions images from multiple or even all forget concepts onto a
340 single base concept. As shown in Table 4, CORE’s one-to-one reconditioning scheme significantly
341 outperforms these strategies. Specifically, unlearning accuracy declines sharply when forget concepts
342 are paired with multiple alternatives (one-to-all or one-to-five) and a random alternative is sampled
343 at each step. Conversely, the model’s stylistic retainability suffers when all forget concepts are
344 reconditioned to just one or a few base concepts.

345 5 Conclusion and Future Directions

346 In this paper, we introduce COncept REconditioning (CORE), a novel and effective method for
347 unlearning in diffusion models. CORE leverages a non-trainable target noise from the pre-trained dif-
348 fusion model to guide both the unlearning and retain objectives, thereby avoiding the pitfalls of using
349 trainable noise predictors or random Gaussian noise targets. Through extensive experiments on the
350 UnlearnCanvas dataset, we demonstrate that CORE consistently outperforms state-of-the-art baseline
351 methods in terms of unlearning effectiveness, retainability, and generalization ability, particularly in
352 challenging tasks involving multiple forget concepts. Moreover, we highlight the importance of a
353 one-to-one concept reconditioning scheme, which proves superior to other anchor selection strategies.
354 There are several promising directions for future research. One key area is improving the efficiency
355 of unlearning, particularly when dealing with a large number of forget concepts. Current methods
356 can still be time-consuming when unlearning many concepts simultaneously. Exploring accelerated
357 unlearning methods while maintaining performance is an exciting avenue. Additionally, future work
358 could investigate the robustness of unlearning methods in dynamic environments, where new concepts
359 might continuously be added to the model, requiring continuous updates without retraining from
360 scratch.

361 **References**

- 362 Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. Multimodal datasets: Misogyny,
363 pornography, and malignant stereotypes.” arxiv, 2021.
- 364 Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers,
365 Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE Symposium*
366 *on Security and Privacy (SP)*, pp. 141–159. IEEE, 2021.
- 367 Brittan. Ai-created images lose u.s. copyrights in test for new technology. *Reuters*, 2023.
- 368 Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In *2015*
369 *IEEE symposium on security and privacy*, pp. 463–480. IEEE, 2015.
- 370 Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr, Borja
371 Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd*
372 *USENIX Security Symposium (USENIX Security 23)*, pp. 5253–5270, 2023.
- 373 CCPA. California consumer privacy act of 2018. https://leginfo.ca.gov/faces/billTextClient.xhtml?bill_id=201720180AB375, 2018. AB-375,
374 Signed into law on June 28, 2018.
375
- 376 Jiaao Chen and Diyi Yang. Unlearn what you want to forget: Efficient unlearning for llms. *arXiv*
377 *preprint arXiv:2310.20150*, 2023.
- 378 Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models
379 in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9):
380 10850–10869, 2023.
- 381 Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale.
382 *arXiv preprint arXiv:2010.11929*, 2020.
- 383 Guangyao Dou, Zheyuan Liu, Qing Lyu, Kaize Ding, and Eric Wong. Avoiding copyright infringe-
384 ment via machine unlearning, 2024. URL <https://arxiv.org/abs/2406.10952>.
- 385 Ronen Eldan and Mark Russinovich. Who’s harry potter? approximate unlearning in llms. *arXiv*
386 *preprint arXiv:2310.02238*, 2023.
- 387 Chongyu Fan, Jiancheng Liu, Yihua Zhang, Dennis Wei, Eric Wong, and Sijia Liu. Salun: Em-
388 powering machine unlearning via gradient-based weight saliency in both image classification and
389 generation. *arXiv preprint arXiv:2310.12508*, 2023.
- 390 Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts
391 from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer*
392 *Vision*, pp. 2426–2436, 2023.
- 393 Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified
394 concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on*
395 *Applications of Computer Vision*, pp. 5111–5120, 2024.
- 396 Chongyang Gao, Lixu Wang, Chenkai Weng, Xiao Wang, and Qi Zhu. Practical unlearning for large
397 language models, 2024. URL <https://arxiv.org/abs/2407.10223>.
- 398 Alvin Heng and Harold Soh. Continual learning for forgetting in deep generative models. 2023.
- 399 Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep
400 generative models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 401 Masayuki Henmi and Shinto Eguchi. A paradox concerning nuisance parameters and projected
402 estimating functions. *Biometrika*, 91(4):929–941, 2004.
- 403 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans
404 trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural*
405 *Information Processing Systems (NeurIPS)*, pp. 6626–6637, 2017.

- 406 Kohtaro Hitomi, Yoshihiko Nishiyama, and Ryo Okui. A puzzling phenomenon in semiparametric
407 estimation problems with infinite-dimensional nuisance parameters. *Econometric Theory*, 24(6):
408 1717–1728, 2008.
- 409 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
410 *neural information processing systems*, 33:6840–6851, 2020.
- 411 Zachary Izzo, Mary Anne Smart, Kamalika Chaudhuri, and James Zou. Approximate data deletion
412 from machine learning models. In *International Conference on Artificial Intelligence and Statistics*,
413 pp. 2008–2016. PMLR, 2021.
- 414 Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and
415 Minjoon Seo. Knowledge unlearning for mitigating privacy risks in language models. *arXiv*
416 *preprint arXiv:2210.01504*, 2022.
- 417 Jiabao Ji, Yujian Liu, Yang Zhang, Gaowen Liu, Ramana Rao Kompella, Sijia Liu, and Shiyu Chang.
418 Reversing the forget-retain objectives: An efficient llm unlearning framework from logit difference,
419 2024. URL <https://arxiv.org/abs/2406.08607>.
- 420 Amirhossein Kazerouni, Ehsan Khodapanah Aghdam, Moein Heidari, Reza Azad, Mohsen Fayyaz,
421 Ilker Hacihaliloglu, and Dorit Merhof. Diffusion models in medical imaging: A comprehensive
422 survey. *Medical Image Analysis*, 88:102846, 2023.
- 423 Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan
424 Zhu. Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF*
425 *International Conference on Computer Vision*, pp. 22691–22702, 2023.
- 426 Agostina J Larrazabal, Nicolás Nieto, Victoria Peterson, Diego H Milone, and Enzo Ferrante. Gender
427 imbalance in medical imaging datasets produces biased classifiers for computer-aided diagnosis.
428 *Proceedings of the National Academy of Sciences*, 117(23):12592–12594, 2020.
- 429 Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D Li,
430 Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, et al. The wmdp benchmark: Measuring
431 and reducing malicious use with unlearning. *arXiv preprint arXiv:2403.03218*, 2024a.
- 432 Senmao Li, Joost van de Weijer, Taihang Hu, Fahad Shahbaz Khan, Qibin Hou, Yaxing Wang, and
433 Jian Yang. Get what you want, not what you don’t: Image content suppression for text-to-image
434 diffusion models. *arXiv preprint arXiv:2402.05375*, 2024b.
- 435 Liu. The world’s smartest artificial intelligence just made its first magazine cover. *Cosmopolitan*,
436 2022.
- 437 Zheyuan Liu, Guangyao Dou, Zhaoxuan Tan, Yijun Tian, and Meng Jiang. Machine unlearning in
438 generative ai: A survey. *arXiv preprint arXiv:2407.20516*, 2024a.
- 439 Zheyuan Liu, Guangyao Dou, Zhaoxuan Tan, Yijun Tian, and Meng Jiang. Towards safer large
440 language models through machine unlearning. *arXiv preprint arXiv:2402.10058*, 2024b.
- 441 Mengyao Lyu, Yuhong Yang, Haiwen Hong, Hui Chen, Xuan Jin, Yuan He, Hui Xue, Jungong Han,
442 and Guiguang Ding. One-dimensional adapter to rule them all: Concepts diffusion models and
443 erasing applications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
444 *Recognition*, pp. 7559–7568, 2024.
- 445 Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C Lipton, and J Zico Kolter. Tofu: A task of
446 fictitious unlearning for llms. *arXiv preprint arXiv:2401.06121*, 2024.
- 447 Alessandro Mantelero. The eu proposal for a general data protection regulation and the roots of the
448 ‘right to be forgotten’. *Computer Law & Security Review*, 29(3):229–235, 2013.
- 449 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual
450 associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- 451 Yisroel Mirsky and Wenke Lee. The creation and detection of deepfakes: A survey. *ACM computing*
452 *surveys (CSUR)*, 54(1):1–41, 2021.

- 453 Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. Memory-
454 based model editing at scale. In *International Conference on Machine Learning*, pp. 15817–15831.
455 PMLR, 2022.
- 456 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
457 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
458 text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 459 Yong-Hyun Park, Sangdoon Yun, Jin-Hwa Kim, Junho Kim, Geonhui Jang, Yonghyun Jeong, Junghyo
460 Jo, and Gayoung Lee. Direct unlearning optimization for robust and safe text-to-image models,
461 2024. URL <https://arxiv.org/abs/2407.21035>.
- 462 Martin Pawelczyk, Seth Neel, and Himabindu Lakkaraju. In-context unlearning: Language models
463 as few shot unlearners. *arXiv preprint arXiv:2310.07579*, 2023.
- 464 Popli. He used ai to publish a children’s book in a weekend. artists are not happy about it. *Time*
465 *Magazine*, 2022.
- 466 Javier Rando, Daniel Paleka, David Lindner, Lennart Heim, and Florian Tramèr. Red-teaming the
467 stable diffusion safety filter. *arXiv preprint arXiv:2210.04610*, 2022.
- 468 James M Robins, Steven D Mark, and Whitney K Newey. Estimating exposure effects by modelling
469 the expectation of exposure conditional on confounders. *Biometrics*, pp. 479–495, 1992.
- 470 Robin Rombach and Patrick Esser. Stablediffusion 2.0. 2022.
- 471 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
472 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-*
473 *ence on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, June 2022.
- 474 Roose. An A.I.-generated picture won an art prize. artists aren’t happy. *The New York Times*, 2022.
- 475 Scenario. Scenario.gg. ai-generated game assets. 2022.
- 476 Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion:
477 Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF*
478 *Conference on Computer Vision and Pattern Recognition*, pp. 22522–22531, 2023.
- 479 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
480 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
481 open large-scale dataset for training next generation image-text models. *Advances in Neural*
482 *Information Processing Systems*, 35:25278–25294, 2022.
- 483 Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember
484 what you want to forget: Algorithms for machine unlearning. *Advances in Neural Information*
485 *Processing Systems*, 34:18075–18086, 2021.
- 486 Shawn Shan, Jenna Cryan, Emily Wenger, Haitao Zheng, Rana Hanocka, and Ben Y Zhao. Glaze:
487 Protecting artists from style mimicry by {Text-to-Image} models. In *32nd USENIX Security*
488 *Symposium (USENIX Security 23)*, pp. 2187–2204, 2023.
- 489 SmithMano. Tutorial: How to remove the safety filter in 5 seconds. *Reddit*, 2023.
- 490 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
491 learning using nonequilibrium thermodynamics. In *International conference on machine learning*,
492 pp. 2256–2265. PMLR, 2015.
- 493 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion
494 art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the*
495 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6048–6058, 2023a.
- 496 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Under-
497 standing and mitigating copying in diffusion models. *Advances in Neural Information Processing*
498 *Systems*, 36:47783–47803, 2023b.

- 499 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
500 *preprint arXiv:2010.02502*, 2020.
- 501 Fangzhou Su, Wenlong Mou, Peng Ding, and Martin J Wainwright. When is the estimated propensity
502 score better? high-dimensional analysis and bias correction. *arXiv preprint arXiv:2303.17102*,
503 2023.
- 504 Anvith Thudi, Gabriel Deza, Varun Chandrasekaran, and Nicolas Papernot. Unrolling sgd: Under-
505 standing factors influencing machine unlearning. In *2022 IEEE 7th European Symposium on*
506 *Security and Privacy (EuroS&P)*, pp. 303–319. IEEE, 2022.
- 507 Bozhong Tian, Xiaozhuan Liang, Siyuan Cheng, Qingbin Liu, Mengru Wang, Dianbo Sui, Xi Chen,
508 Huajun Chen, and Ningyu Zhang. To forget or not? towards practical knowledge unlearning for
509 large language models, 2024. URL <https://arxiv.org/abs/2407.01920>.
- 510 Paul Voigt and Axel Von dem Bussche. The eu general data protection regulation (gdpr). *A Practical*
511 *Guide, 1st Ed., Cham: Springer International Publishing*, 10(3152676):10–5555, 2017.
- 512 Nikhil Vyas, Sham M Kakade, and Boaz Barak. On provable copyright protection for generative
513 models. In *International Conference on Machine Learning*, pp. 35277–35299. PMLR, 2023.
- 514 Jing Wu and Mehrtash Harandi. Scissorhands: Scrub data influence via connection sensitivity in
515 networks. *arXiv preprint arXiv:2401.06187*, 2024.
- 516 Jing Wu, Trung Le, Munawar Hayat, and Mehrtash Harandi. Erasediff: Erasing data influence in
517 diffusion models. *arXiv preprint arXiv:2401.05779*, 2024.
- 518 Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geometric
519 diffusion model for molecular conformation generation. *arXiv preprint arXiv:2203.02923*, 2022.
- 520 Ling Yang and Shenda Hong. Unsupervised time-series representation learning with iterative bilinear
521 temporal-spectral fusion. In *International conference on machine learning*, pp. 25038–25054.
522 PMLR, 2022.
- 523 Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. *arXiv preprint*
524 *arXiv:2310.10683*, 2023.
- 525 Hongbang Yuan, Zhuoran Jin, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. Towards ro-
526 bust knowledge unlearning: An adversarial framework for assessing and improving unlearning
527 robustness in large language models, 2024. URL <https://arxiv.org/abs/2408.10682>.
- 528 Gong Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not:
529 Learning to forget in text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference*
530 *on Computer Vision and Pattern Recognition*, pp. 1755–1764, 2024a.
- 531 Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From catastrophic
532 collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*, 2024b.
- 533 Yihua Zhang, Yimeng Zhang, Yuguang Yao, Jinghan Jia, Jiancheng Liu, Xiaoming Liu, and Sijia Liu.
534 Unlearncanvas: A stylized image dataset to benchmark machine unlearning for diffusion models.
535 *arXiv preprint arXiv:2402.11846*, 2024c.
- 536 Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and
537 Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate
538 unsafe images ... for now, 2024d. URL <https://arxiv.org/abs/2310.11868>.
- 539 Zhexin Zhang, Junxiao Yang, Pei Ke, Shiyao Cui, Chujie Zheng, Hongning Wang, and Minlie Huang.
540 Safe unlearning: A surprisingly effective and generalizable solution to defend against jailbreak
541 attacks, 2024e. URL <https://arxiv.org/abs/2407.02855>.

542 A Related Works

543 **Malicious Behavior of Diffusion Models.** Diffusion models have demonstrated impressive ca-
544 pabilities in generating high-quality, efficient text-to-image outputs (Ho et al., 2020; Song et al.,
545 2020; Rombach et al., 2022). However, these large-scale trained models can pose significant privacy
546 and ethical risks. They are capable of memorizing private images and reproducing objectionable
547 content, such as pornography, private personal photos, malign stereotypes, gender and race bi-
548 ases (Schramowski et al., 2023; Larrazabal et al., 2020; Carlini et al., 2023; Somepalli et al., 2023a;
549 Rando et al., 2022). This mainly stems from the contaminated data sources which involves problem-
550 atic image-text pairs (Birhane et al., 2021). Furthermore, diffusion models can cause potential issues
551 about copyright infringement by mimicking, or even replicating the styles of some specific artistic
552 and their copyrighted work (Shan et al., 2023). Reports showed that AI-generated arts can sometimes
553 be published commercially (Liu, 2022; Popli, 2022; Scenario, 2022) and even awarded prizes (Roose,
554 2022), raising more serious social concerns about intellectual property violations (Brittan, 2023;
555 Somepalli et al., 2023b; Shan et al., 2023).

556 **Diffusion Model Unlearning.** The goal of unlearning diffusion models is to eliminate unwanted
557 concepts and their influence on model outputs. Directly retraining a model to remove such concepts is
558 highly resource-intensive and thus inefficient for large diffusion models (Nichol et al., 2021; Rombach
559 et al., 2022; Schuhmann et al., 2022). Recent research has explored more efficient unlearning
560 techniques. One approach focuses on inference-time methods, which attempt to filter or steer the
561 model away from undesirable outputs during generation (Rando et al., 2022; Schramowski et al.,
562 2023). However, these methods are often limited in effectiveness and can be bypassed, particularly in
563 open-source models (SmithMano, 2023). A more robust alternative involves fine-tuning the model’s
564 parameters to actively remove undesirable concepts from its learned representations (Zhang et al.,
565 2024a; Li et al., 2024b; Lyu et al., 2024; Heng & Soh, 2023; Vyas et al., 2023; Gandikota et al., 2024).
566 Some methods are similar to ours: Gandikota et al. (2023); Wu et al. (2024) match the denoising
567 network of correct images given a target concept to another distribution. Fan et al. (2023) additionally
568 adds a saliency map to fine-tune only a small fraction of parameters. Heng & Soh (2024) does gradient
569 ascent on the training loss of diffusion models. Kumari et al. (2023) minimizes the distribution
570 mismatch between the target concept and another anchor concept. We will discuss the difference
571 between our algorithm and theirs in more detail in Section 3.2. Effective though, achieving robust
572 unlearning on complex tasks still remains challenging (Zhang et al., 2024c,d). For a comprehensive
573 review of unlearning techniques in generative models, see Liu et al. (2024a).

574 **Machine Unlearning.** Machine unlearning has been extensively explored within classification
575 tasks (Cao & Yang, 2015; Bourtole et al., 2021; Sekhari et al., 2021; Izzo et al., 2021; Thudi et al.,
576 2022) and is now being applied to large generative models. One popular class of unlearning methods
577 stems from Gradient Ascent(GA) (Jang et al., 2022; Yao et al., 2023; Chen & Yang, 2023; Zhang et al.,
578 2024b). More methods include preference optimization (Zhang et al., 2024b; Maini et al., 2024; Park
579 et al., 2024), model-editing (Meng et al., 2022; Mitchell et al., 2022; Eldan & Russinovich, 2023),
580 knowledge negation (Liu et al., 2024b), representation control (Li et al., 2024a), logits difference
581 method (Ji et al., 2024), random labeling, saliency map (Dou et al., 2024; Tian et al., 2024), and
582 in-context unlearning approaches (Pawelczyk et al., 2023), etc. Some other methods are developed
583 for adversarial unlearning or sequential unlearning tasks (Zhang et al., 2024e; Yuan et al., 2024; Gao
584 et al., 2024). These unlearning methods for language models are orthogonal to our proposed method
585 for unlearning diffusion models.

586 B Experiment Details

587 **Hyperparameter.** All experiments are done using one 80GB NVIDIA A100 GPU. We use an open-
588 sourced Stable Diffusion v-1.5 for all experiments (Rombach et al., 2022), which is first fine-tuned
589 on all data in UnlearnCanvas before any unlearning process, and the fine-tuned model is provided
590 by Zhang et al. (2024c). As suggested in prior works (Gandikota et al., 2023; Zhang et al., 2024c),
591 we only fine-tune the cross-attention in U-Nets in the Stable Diffusion and freeze all other parameters
592 when doing unlearning. Following Zhang et al. (2024c), we use the first three images for each style
593 and object for training. For CORE, we run 25 epochs in Forget01 and Forget06, and 100 epochs in
594 Forget25. In testing the generalization ability of unlearning styles, where the testing and training
595 objects are distinct, we double the epochs in Forget06. We use Adam with a constant learning
596 rate of 1×10^{-5} in CORE, and the batch size is set to 4. We set $\alpha = 1.0$ in equation (4). The
597 hyperparameters used for training the baseline methods are described in Appendix C.

598 **Scaled FID Values.** Scaled FID (**SFID**) is a modified version of Fréchet Inception Distance
599 (FID) (Heusel et al., 2017), which ranges from zero to infinity and measures the quality of generated
600 images. A lower FID value indicates a higher generation quality. To measure the overall performance
601 of unlearning algorithms, we convert the original FID value into Scaled FID value, which ranges
602 from 0 to 100 and increases when the generation quality grows. We compute the original FID value
603 for the base model and the unlearned model, denoted as \mathbf{FID}_0 and \mathbf{FID}_M , respectively. SFID is then
604 defined as

$$\mathbf{SFID}_M = \min \left\{ 100 \times \frac{\mathbf{FID}_0}{\mathbf{FID}_M}, 100 \right\} \quad (7)$$

605 A model with better retainability tends to have higher SFID values. In our experiments, we compute
606 SFID values on the retain set.

607 C Baseline Methods Overview

608 In this section, we introduce baseline methods, discuss how they relate to our proposed approach, and
 609 describe their training procedures. For the most part, the training setup for these baseline methods
 610 follows Zhang et al. (2024c). We set the alternative concept as one common base concept (one base
 611 style) in Forget01. For each step, we randomly sample one alternative concept from the retain set in
 612 Forget06. In Forget25, we create a bijection from the 25 concepts in the forget set and the other 25
 613 concepts in the retain set. In Forget25, we have also tried to pick a random alternative concept at
 614 each step, but this worsens the performance for all baselines by a large margin.

615 **ESD (Gandikota et al., 2023).** ESD is the first method that offers both efficiency and effectiveness
 616 in unlearning for diffusion models. It utilizes a more complex unlearning objective without incorpo-
 617 rating a retain objective. As a result, ESD’s retainability is generally outperformed by other methods.
 618 The objective function for ESD is defined as follows:

$$\mathcal{L}_{\text{ESD}}(\theta) := \mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_{f,t}} \left\| \varepsilon_{\theta}(\mathbf{x}_t, p_f) - \left(\varepsilon_{\theta^*}(\mathbf{x}_t, p_0) - \eta(\varepsilon_{\theta^*}(\mathbf{x}_t, p_f) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_0)) \right) \right\|_2^2, \quad (8)$$

619 where (\mathbf{x}_0, p_f) are sampled from the forget set, t is uniformly sampled from $[0, T]$, ε_{θ} and ε_{θ^*} are
 620 the current and pre-trained noise predictors in diffusion models. Here, p_0 is a base concept, which
 621 can be an empty string (Gandikota et al., 2023) or a base style in UnlearnCanvas. In our experiments,
 622 according to Gandikota et al. (2023), we set $\eta = 1.0$, batch size to 1, and the learning rate to 1×10^{-5} ,
 623 and we run 1000 gradient steps.

624 Although the objective in ESD seems to be very different from our framework of concept erasing, we
 625 can still fit it into our framework in Section 3.2 via proper decomposition. Namely,

$$\begin{aligned} \mathcal{L}_{\text{ESD}}(\theta) &:= \mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_{f,t}} \left\| \varepsilon_{\theta}(\mathbf{x}_t, p_f) - \left(\varepsilon_{\theta^*}(\mathbf{x}_t, p_0) - \eta(\varepsilon_{\theta^*}(\mathbf{x}_t, p_f) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_0)) \right) \right\|_2^2 \\ &= \mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_{f,t}} \left\| \varepsilon_{\theta}(\mathbf{x}_t, p_f) - (1 + \eta)\varepsilon_{\theta^*}(\mathbf{x}_t, p_0) + \eta\varepsilon_{\theta^*}(\mathbf{x}_t, p_f) \right\|_2^2 \\ &= \mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_{f,t}} \left(\varepsilon_{\theta}(\mathbf{x}_t, p_f)^2 + (1 + \eta)^2 \cdot \varepsilon_{\theta^*}(\mathbf{x}_t, p_0)^2 + \eta^2 \cdot \varepsilon_{\theta^*}(\mathbf{x}_t, p_f)^2 \right. \\ &\quad \left. + 2\eta \cdot \varepsilon_{\theta}(\mathbf{x}_t, p_f) \cdot \varepsilon_{\theta^*}(\mathbf{x}_t, p_f) - 2(1 + \eta) \cdot \varepsilon_{\theta}(\mathbf{x}_t, p_f) \cdot \varepsilon_{\theta^*}(\mathbf{x}_t, p_0) \right. \\ &\quad \left. - 2\eta(1 + \eta) \cdot \varepsilon_{\theta^*}(\mathbf{x}_t, p_0) \cdot \varepsilon_{\theta^*}(\mathbf{x}_t, p_f) \right) \\ &= \mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_{f,t}} \left((1 + \eta) \cdot (\varepsilon_{\theta}(\mathbf{x}_t, p_f) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_0))^2 - \eta \cdot (\varepsilon_{\theta}(\mathbf{x}_t, p_f) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_f))^2 \right. \\ &\quad \left. + \eta(1 + \eta) \cdot (\varepsilon_{\theta^*}(\mathbf{x}_t, p_0) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_f))^2 \right) \\ &= \underbrace{(1 + \eta) \mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_{f,t}} (\varepsilon_{\theta}(\mathbf{x}_t, p_f) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_0))^2}_{(a)} \\ &\quad - \underbrace{\eta \mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_{f,t}} (\varepsilon_{\theta}(\mathbf{x}_t, p_f) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_f))^2}_{(b)} \\ &\quad + \underbrace{\eta(1 + \eta) \mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_{f,t}} (\varepsilon_{\theta^*}(\mathbf{x}_t, p_0) - \varepsilon_{\theta^*}(\mathbf{x}_t, p_f))^2}_{(c)}. \end{aligned}$$

626 Since term (c) in the last line is a constant independent of θ , we can omit it in the loss function. The
 627 remaining two terms (a) and (b) can both fit into the Concept Erasing framework (see equation 5).
 628 Term (a) is equivalent to choosing $\lambda = (1 + \eta)$ and $\mathbf{y}_{\text{CE}} = \varepsilon_{\theta^*}(\mathbf{x}_t, p_0)$, while term (b) is equivalent
 629 to choosing $\lambda = -\eta$ and $\mathbf{y}_{\text{CE}} = \varepsilon_{\theta^*}(\mathbf{x}_t, p_f)$.

630 **SalUn (Fan et al., 2023).** Saliency Unlearning (SalUn) introduces a saliency mask to the diffusion
 631 model parameters before unlearning. This mask, based on the absolute gradient scale for the forget
 632 concept, identifies the most important parameter subsets for unlearning targeted concepts, enabling
 633 efficient unlearning that edits only a small portion of the model. The loss function for SalUn is given

634 by:

$$\mathcal{L}_{\text{SalUn}}(\theta) := \underbrace{\mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_f, t, p_r \neq p_f} \|\varepsilon_\theta(\mathbf{x}_t, p_f) - \varepsilon_\theta(\mathbf{x}_t, p_r)\|_2^2}_{\text{unlearn objective}} + \underbrace{\beta \cdot \mathbb{E}_{(\mathbf{x}_0, p_r) \sim \mathcal{D}_r, t, \varepsilon} \|\varepsilon - \varepsilon_\theta(\mathbf{x}_t, p_r)\|_2^2}_{\text{retain objective}}, \quad (9)$$

635 where ε is the standard Gaussian random vector used to generate \mathbf{x}_t , and p_f and p_r are forget concepts
 636 and retain concepts, respectively. t is sampled uniformly from $[0, T]$. In contrast to CORE, which
 637 uses ε_{θ^*} as the target for the retain objectives, SalUn uses the Gaussian random vector ε . Their
 638 unlearn objective can fit in the framework in equation (5) with a trainable network as the target noise.
 639 This can lead to target degradation during the unlearning process, especially when multiple concepts
 640 need to be unlearned. Following Fan et al. (2023) and Zhang et al. (2024c), we take $\beta = 1.0$. We use
 641 a learning rate of 1×10^{-5} and a batch size of 4. We run 10 epochs in Forget01 and 100 epochs in
 642 Forget06 and Forget25.

643 **EDiff (Wu et al., 2024).** EraseDiff (EDiff) formulates the objective as follows:

$$\mathcal{L}_{\text{EDiff}}(\theta) := \underbrace{\mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_f, t, \varepsilon_f} \|\varepsilon_\theta(\mathbf{x}_t, p_f) - \varepsilon_f\|_2^2}_{\text{unlearn objective}} + \underbrace{\beta \cdot \mathbb{E}_{(\mathbf{x}_0, p_r) \sim \mathcal{D}_r, t, \varepsilon} \|\varepsilon - \varepsilon_\theta(\mathbf{x}_t, p_r)\|_2^2}_{\text{retain objective}}. \quad (10)$$

644 The retain objective is similar to that in SalUn, but the unlearn objective differs. Here, ε_f is a uni-
 645 formly distributed random vector, which serves as the target noise. This unlearn objective aligns with
 646 the concept erasing framework (equation 5), where \mathbf{y}_{CE} is uniformly distributed. EraseDiff simplifies
 647 the diffusion process by solving it as a first-order optimization problem, reducing computational
 648 complexity. In our experiments, we use a batch size of 4 and a learning rate of 5×10^{-5} . We run 5
 649 epochs in Forget01 and 50 epochs in Forget06 and Forget25.

650 **CA (Kumari et al., 2023).** Concept Ablation (CA) matches the image distribution from the forget
 651 set to an anchor concept. They design two objective functions: a model-based one and a noise-based
 652 one. The model-based CA objective is defined as

$$\mathcal{L}_{\text{CA-model}}(\theta) := \underbrace{\mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_f, t} [\omega_t \|\varepsilon_\theta(\mathbf{x}_t, p_f) - \varepsilon_\theta(\mathbf{x}_t, p_0) \cdot \text{sg}(\cdot)\|_2]}_{\text{unlearn objective}} + \underbrace{\lambda \cdot \mathbb{E}_{(\mathbf{x}_0, p_r) \sim \mathcal{D}_r, t, \varepsilon} \|\varepsilon - \varepsilon_\theta(\mathbf{x}_t, p_r)\|_2^2}_{\text{retain objective}}. \quad (11)$$

653 Here, ω_t is a time-dependent weight applied to the loss, p_0 is a fixed base concept from the retain set,
 654 and $\text{sg}(\cdot)$ denotes the stop-gradient operator. The noise-based objective is defined as

$$\mathcal{L}_{\text{CA-noise}}(\theta) := \underbrace{\mathbb{E}_{(\mathbf{x}_0, p_f) \sim \mathcal{D}_f, t, \varepsilon} [\omega_t \|\varepsilon_\theta(\mathbf{x}_t, p_f) - \varepsilon\|_2]}_{\text{unlearn objective}} + \underbrace{\lambda \cdot \mathbb{E}_{(\mathbf{x}_0, p_r) \sim \mathcal{D}_r, t, \varepsilon} \|\varepsilon - \varepsilon_\theta(\mathbf{x}_t, p_r)\|_2^2}_{\text{retain objective}}. \quad (12)$$

655 In both objectives, ε is the standard Gaussian random vector used to generate \mathbf{x}_t . In our experiments,
 656 we use a batch size of 4 and a learning rate of 1.6×10^{-5} . We run 200 gradient steps in Forget01 and
 657 100 epochs in Forget06 and Forget25.



Figure 3: Additional generated images from the unlearned model. The first column is generated by the fine-tuned Stable Diffusion model before any unlearning. Other columns are generated by the model unlearned by our proposed method and five baseline methods.

658 **D More results**

659 In this section, we present more images generated from our experiments on UnlearnCanvas in Figure 3.

660