

CONCEPTPRUNE: CONCEPT EDITING IN DIFFUSION MODELS VIA SKILLED NEURON PRUNING

Anonymous authors

Paper under double-blind review

ABSTRACT

While large-scale text-to-image diffusion models have demonstrated impressive image-generation capabilities, there are significant concerns about their potential misuse for generating unsafe content, violating copyright, and perpetuating societal biases. Recently, the text-to-image generation community has begun addressing these concerns by editing or unlearning undesired concepts from pre-trained models. However, these methods often involve data-intensive and inefficient fine-tuning or utilize various forms of token remapping, rendering them susceptible to adversarial jailbreaks. In this paper, we present a simple and effective training-free approach, *ConceptPrune*, wherein we first identify critical regions within pre-trained models responsible for generating undesirable concepts, thereby facilitating straightforward concept unlearning via weight pruning. Experiments across a range of concepts including artistic styles, nudity, and object erasure demonstrate that target concepts can be efficiently erased by pruning a tiny fraction, approximately 0.12% of total weights, enabling multi-concept erasure and robustness against various white-box and black-box adversarial attacks.

1 INTRODUCTION

In recent years, text-to-image generation has witnessed significant advances driven by the development and adoption of diffusion models (DMs) [Ho et al., 2020; Rombach et al., 2021; Ruiz et al., 2022; Saharia et al., 2022; Nichol et al., 2021; Zhang et al., 2023c; Luo et al., 2023; Podell et al., 2023] across industries and real-world scenarios. However, this swift advancement presents a substantial risk. Diffusion models can threaten artists’ livelihoods through style replication [et al v. Stability AI Ltd. et al., 2023], generate convincing deepfakes and NSFW content [Review, 2023; Forensics, 2024], and perpetuate societal biases [Luccioni et al., 2023]. The risks associated with large-scale text-to-image models arise from billion-sized web-scraped datasets used in training, comprising public datasets like LAION [Schuhmann et al., 2022], COYO [Byeon et al., 2022], and CC12M [Changpinyo et al., 2021], that often lack human-level quality assurance. A simplistic and naive solution to mitigate these risks involves fine-tuning the model on datasets without this undesired content; however, this approach can prove to be highly compute-expensive.

Several efforts addressing the risks of diffusion models have been made from the perspective of Concept Editing [Kumari et al., 2023; Gandikota et al., 2023a;b; Zhang et al., 2023a; Orgad et al., 2023] and Model Unlearning (MU) [Heng & Soh, 2023; Zhao et al., 2024; Liu et al., 2024; Wu et al., 2024; Fan et al., 2023], both aimed at eliminating undesired prompts, albeit with differing objectives. Concept editing methods seek to eliminate undesired prompts by aligning latent representations of the target concept with a concept to be retained, via methods such as maximizing similarity [Kumari et al., 2023; Gandikota et al., 2023a] and token remapping [Zhang et al., 2023a; Gandikota et al., 2023b]. Conversely, Model Unlearning formulates an objective that penalizes forgetting desired concepts while promoting the elimination of undesired ones, but this requires expensive computations and fine-tuning. Moreover, as most concept editing approaches rely on some form of token blacklisting or re-steering [Zhang et al., 2023a], adversarial attacks based on textual inversion [Zhang et al., 2023d; Pham et al., 2023; Yang et al., 2023; Tsai et al., 2024] have demonstrated the ability to circumvent concept erasure methods [Gandikota et al., 2023a;b; Zhang et al., 2023a] that were previously believed to be robust with a near-perfect success rate.

054 In this paper, we introduce *ConceptPrune*, an entirely training-free method for concept editing
055 that, for the first time, tackles knowledge editing in diffusion models through the lens of pruning.
056 Leveraging recently introduced pruning heuristics [Sun et al., 2024], we identify regions or neurons in
057 feed-forward layers of diffusion models that strongly activate in the presence of a concept, and denote
058 them as *skilled neurons*. Subsequently, concept removal can be achieved by simply pruning or *zeroing*
059 out these skilled regions. We demonstrate that ConceptPrune provides a rapid, efficient, and unified
060 solution for erasing undesired concepts, including various artist styles, nudity, undesired objects, and
061 gender biases. Notably, it maintains the outstanding image-generation prowess of pre-trained models
062 while remaining resilient to adversarial attacks.

063 064 2 RELATED WORK

065
066
067 **Concept Erasure in Diffusion Models:** Concept erasure has gained significant attention and has
068 rapidly emerged as a pivotal area of research in diffusion models. Recent concept erasure methods
069 can be broadly categorized into two main areas: *Model Unlearning* and *Concept Editing*.

070 Model Unlearning methods [Heng & Soh, 2023; Wu & Harandi, 2024; Fan et al., 2023; Zhang
071 et al., 2024] typically require extensive training to forget a target concept while preserving unrelated
072 ones. While these methods have shown remarkable efficacy in unlearning multiple concepts, they are
073 usually computationally expensive, especially for large-scale models.

074 Concept Editing [Gandikota et al., 2023a;b; Kumari et al., 2023; Zhang et al., 2023a; Huang et al.,
075 2023; Lu et al., 2024; Lyu et al., 2023] focuses on modifications to specific parts of the model. These
076 edits ensure that the denoised output for the target concept aligns with clean, desired concepts. The
077 training costs associated with Concept Editing can be mitigated by strategies such as tuning only
078 cross-attention weights [Gandikota et al., 2023a; Kumari et al., 2023; Zhang et al., 2023a; Huang
079 et al., 2023; Lu et al., 2024], solving closed-form objectives to update attention parameters [Gandikota
080 et al., 2023b; Lu et al., 2024; Orgad et al., 2023], or parameter-efficient adaptation like LORA [Hu
081 et al., 2022] to edit the model [Lu et al., 2024; Lyu et al., 2023].

082 While the aforementioned methods are highly effective, deploying current state-of-the-art concept
083 erasure techniques in real-world scenarios poses significant challenges, particularly in online envi-
084 ronments with computational constraints where harmful concepts can emerge dynamically. This is
085 because these methods struggle to meet the following requirements for real-world applications: (1)
086 *training-free concept erasure*, eliminating concepts without the need for backpropagation through the
087 entire model, or (2) *lightweight or fast concept erasure*, allowing concepts to be removed quickly and
088 efficiently with minimal compute.

089 Most concept-erasure methods rely on extensive fine-tuning and are therefore not training-free,
090 however some training-based approaches like UCE [Gandikota et al., 2023b], SPM [Lyu et al., 2023],
091 MACE [Lyu et al., 2023], and FMN [Zhang et al., 2023a] are notably *lightweight* and suitable for the
092 real-world setting. For instance, FMN erases concepts in 30 seconds and UCE in about 2 minutes,
093 while the rapid fine-tuning of LoRA parameters makes SPM and MACE ideal for real-time online
094 erasure. In Table 1, we present a comprehensive summary of related works, categorizing them based
095 on whether they are training-free and lightweight for an online setting.

096 Our proposed solution, ConceptPrune, excels on both fronts by introducing a training-free, pruning-
097 based approach that eliminates harmful concepts without updating any parameters. Instead, it
098 identifies and targets the neurons responsible for generating these concepts enabling efficient concept
099 erasure with significantly reduced computational requirements.

100 **Language model skilled neuron identification:** Previous works [Wang et al., 2022; Suau et al.,
101 2020; Durrani et al., 2023; Dalvi et al., 2018; Durrani et al., 2020; Antverg & Belinkov, 2022] present
102 strong evidence that activation of specific neurons in feed-forward networks in transformers show
103 high correlation with task labels, with perturbations to these neurons impacting task performance.
104 Modular components within pre-trained transformers were identified by leveraging the inherent
105 sparsity in neurons, as shown in [Zhang et al., 2022]. Further, [Zhang et al., 2023e] demonstrates
106 that these modules are specialized in distinct functions. In this work, we aim to identify neurons
107 accountable for generating undesired concepts in diffusion models — a pursuit hitherto unexplored in
this domain. Unlike language models, identifying neurons in diffusion models is complicated due to

Method	Training-free	Parameters Trained	Lightweight Erasure
CA [Kumari et al., 2023]	×	Full denoiser	×
SA [Heng & Soh, 2023]	×	Full denoiser	×
SH [Wu & Harandi, 2024]	×	Full denoiser	×
AdvUnlearn [Zhang et al., 2024]	×	Full denoiser	×
SalUn [Fan et al., 2023]	×	Full denoiser	×
ESD [Gandikota et al., 2023a]	×	Cross Attention	×
Receler [Huang et al., 2023]	×	Cross Attention	✓
FMN [Zhang et al., 2023a]	×	Cross Attention	✓
SPM [Lyu et al., 2023]	×	LORA	✓
MACE [Lu et al., 2024]	×	Cross Attention + LORA	✓
UCE [Gandikota et al., 2023b]	✓	Cross Attention	✓
Ours (ConceptPrune)	✓	<i>None</i>	✓

Table 1: Summary of recent Concept Erasure baselines. ConceptPrune is a training-free approach that enables rapid pruning of the model to eliminate a new target concept without the need for extensive re-training.

the intricate aggregation of neurons across multiple denoising time steps and the model’s sensitivity to the output of previous time steps.

Language model pruning: Network pruning [LeCun et al., 1989; Liu et al., 2019; Han et al., 2015; Frankle & Carbin, 2019; Blalock et al., 2020] aims to reduce model size either by eliminating parameters and substructures from networks [Li et al., 2017; Frantar & Alistarh, 2023] or by masking parameters guided by a score function [Frantar & Alistarh, 2023; Frantar et al., 2023; Sun et al., 2024; Lee et al., 2019]. This study primarily focuses on the latter approach. Exploration of diffusion model pruning is limited, although one study [Fang et al., 2023] introduces structural pruning by accumulating gradient-based importance scores across a chosen subset of denoising time steps. A study [Wei et al., 2024] explores safety-aligned LLMs that inhibit harmful prompts by leveraging pruning heuristics to identify regions responsible for denying harmful responses. In contrast, we use pruning heuristics to locate critical weight regions responsible for unsafe behaviors in pre-trained models and permanently unlearn them through pruning.

3 PRELIMINARIES

(Latent) diffusion models: Diffusion models (DMs) [Ho et al., 2020; Song et al., 2021] are essentially image denoisers that learn to reverse a forward Markov process in which noise is added into input images for multiple time steps $t \in [0, T]$. During training, given a real image \mathbf{x}_0 , a noisy image \mathbf{x}_t at time t is obtained by $\sqrt{a_t}\mathbf{x}_0 + \sqrt{1 - a_t}\epsilon$, where $\epsilon \sim \mathcal{N}(0, I)$ and a_t is a gradually decaying parameter. Then, the denoiser learns to predict the noise added for obtaining \mathbf{x}_t , such that \mathbf{x}_0 can be reconstructed back by deducting predicted noise from \mathbf{x}_t .

Latent diffusion models (LDMs) [Rombach et al., 2022b; Zhang et al., 2023b] are widely used as the first choice of DMs as they accelerate the above process by operating in a latent space \mathbf{z} , of input \mathbf{x} . Thus, a LDM consists of a latent embedding denoiser $f_\theta(\cdot)$, which is trained to predict the added noise by stochastically minimizing the objective $\mathcal{L}(\mathbf{z}, p) = \mathbb{E}_{\epsilon, \mathbf{x}, p, t} [\|\epsilon - f_\theta(\mathbf{z}_t, p, t)\|]$. Given a text prompt p , an encoder which extracts \mathbf{z}_0 from \mathbf{x}_0 and a decoder which maps the denoised $\hat{\mathbf{z}}_0$ to the pixel space. To synthesize an image during inference based on text prompt p , one first samples a noisy embedding \mathbf{z}_T which is iteratively denoised for T time steps until $\hat{\mathbf{z}}_0$ for generating the final image is obtained. Normally, the encoder and decoder are obtained from a frozen pre-trained autoencoder.

4 CONCEPTPRUNE: A TRAINING-FREE CONCEPT EDITING FRAMEWORK

Motivation: Concept editing methods aim to eliminate the undesired concept from a pretrained DM. Inspired by the observation that concepts can activate specific neurons in a neural network [Mahendran & Vedaldi, 2015; Wang et al., 2022], we ask the question: *Can we remove an undesired concept from a pre-trained DM by simply finding neurons specific to this concept, and pruning them?* The answer is *yes*. We show that neurons in LDMs often specialise to specific concepts, and

that pruning these neurons can be used to permanently eliminate undesired concepts from image generation.

4.1 FEED FORWARD NETWORKS (FFNS) IN LATENT DIFFUSION MODELS

We focus on a pre-trained LDM, i.e. Stable Diffusion [Rombach et al., 2021], characterized by a UNet [Ronneberger et al., 2015] denoted as f_θ . The UNet architecture incorporates two ResNet blocks that sandwich two transformer blocks with self-attention between latent representations, cross-attention for the transfer of information from conditional inputs to latent representations, and a Feed-forward Network (FFN) with GEGLU activation [Shazeer, 2020]. Prior research in concept editing, such as [Gandikota et al., 2023a] and [Zhang et al., 2023a], primarily examines cross-attention or self-attention visualizations to detect concept presence or generation. Diverging from this approach and drawing inspiration from NLP skill discovery [Suau et al., 2020; Wang et al., 2022; Zhang et al., 2023e; Durrani et al., 2020; Dalvi et al., 2018], our focus lies on neurons within the Feed-forward networks.

We begin by denoting the input to the FFN layer l at time step t for text prompt p by $\mathbf{z}_t^l(p) \in \mathbb{R}^{d \times N}$, where N is the number of latent tokens and corresponding output by $\mathbf{z}_t^{l+1}(p) \in \mathbb{R}^{d \times N}$. FFN in Stable Diffusion consists of GEGLU activation [Shazeer, 2020] which operates as shown in Equation 1.

$$\begin{aligned} \mathbf{h}_t^l(p) &= \sigma(\mathbf{W}^{l,1} \cdot \mathbf{z}_t^l(p)) \\ \mathbf{z}_t^{l+1}(p) &= \mathbf{W}^{l,2} \cdot \mathbf{h}_t^l(p) \end{aligned} \quad (1)$$

where, $\mathbf{W}^{l,1} \in \mathbb{R}^{d' \times d}$, $\mathbf{W}^{l,2} \in \mathbb{R}^{d \times d'}$ are weight matrices in the first and second linear layers, bias terms are omitted for simplicity and $\sigma(\cdot)$ is GEGLU activation [Hendrycks & Gimpel, 2023]. In our work, we regard $\mathbf{W}^{l,2}[i, :]$ the i -th row and $\mathbf{W}^{l,2}[i, j]$ the element in i -th row and j -th column of matrix $\mathbf{W}^{l,2}$.

4.2 PRUNING STRATEGY: WANDA

We start with recapping the pruning method Wanda [Sun et al., 2024] for the large language models (LLMs), and its adaptation to diffusion models. We denote the weights of linear layer by $\mathbf{W} \in \mathbb{R}^{d_{out} \times d_{in}}$ and input $\mathbf{X} \in \mathbb{R}^{B \times d_{in}}$, where B is the number of data points, i.e. the number of prompts in this paper. Unlike magnitude-based pruning, which considers the weights' magnitude alone, the concept behind the Wanda score is to estimate the combined effect of weights and the magnitude of features on neuron activations. Therefore, the importance of each weight is calculated as an element-wise product of its magnitude and the corresponding input feature-dimension-wise ℓ_2 norm as shown in Equation 2

$$\mathbf{S}(\mathbf{W}, \mathbf{X}) = |\mathbf{W}| \odot (\mathbf{1}^{d_{out}} \cdot \|\mathbf{X}\|_2) \in \mathbb{R}^{d_{out} \times d_{in}}. \quad (2)$$

Here $|\cdot|$ to denote the absolute value operator, $\|\mathbf{X}\|_2$ computes the ℓ_2 norm of each column of \mathbf{X} and results in a d_{in} dimensional vector, and \odot represents element-wise matrix multiplication. Specifically, Eq 2 broadcasts $\|\mathbf{X}\|_2$ across different rows of \mathbf{W} for computing the element-wise product in each row. For each row of \mathbf{W} , represented by $\mathbf{W}_{i,:}$ with corresponding Wanda score $\mathbf{S}(\mathbf{W}, \mathbf{X})_{i,:}$, the bottom- $k\%$ weights with the lowest scores are zeroed out [Sun et al., 2024]. This process effectively induces sparsity in each row of the weights \mathbf{W} by eliminating the bottom- $k\%$ of the weights, as a row is connected to a single activation in the output of a linear layer as a *per-output basis* [Sun et al., 2024]. Elements of the weight matrix \mathbf{W} are often referred to as *weight neurons*, which are different from neurons corresponding to the output of a layer. After pruning the least important weight neurons in a layer, subsequent layers in the model receive updated input activations. Wanda does not require any costly weight update since it solely relies on a calibration set to compute the feature norm matrix, which can be obtained with just a single forward pass through the model. The following will discuss how we use Wanda to prune each row's top- $k\%$ weight neurons for eliminating a concept.

4.3 IDENTIFYING SKILLED NEURONS IN LATENT DIFFUSION MODELS

Target and reference concept prompts: We first define two sets of calibration prompts $\mathcal{P}^* = \{p_1^*, p_2^*, \dots, p_M^*\}$ and $\mathcal{P} = \{p_1, p_2, \dots, p_M\}$ using M objects that can be generated by the model in target and reference concepts, respectively. Here, p_i^* and p_i represent prompts with the target and

reference concepts, respectively. Objects represent common categories, including ‘cat’, ‘dog’, etc. To eradicate the target concept, e.g., "Van Gogh" painting style, we formulate a p_i^* as ‘a <object> in Van Gogh style’ and a p_i as ‘a <object>’.

Importance score for FFN weights at time t : We begin by collecting the neuron activations described in Eq 1, corresponding to the sets of target concept and reference prompts, and shape them into matrices denoted by $\mathbf{H}_t^l(\mathcal{P}^*) = [\mathbf{h}_t^l(p_1^*)^T, \mathbf{h}_t^l(p_2^*)^T, \dots, \mathbf{h}_t^l(p_M^*)^T]$ and $\mathbf{H}_t^l(\mathcal{P}) = [\mathbf{h}_t^l(p_1)^T, \mathbf{h}_t^l(p_2)^T, \dots, \mathbf{h}_t^l(p_M)^T]$ such that $\mathbf{H}_t^l(\mathcal{P}^*), \mathbf{H}_t^l(\mathcal{P}) \in \mathbf{R}^{(M \times N) \times d'}$. Note that this process only requires one forward pass for per prompt.

After collecting both sets of neuron activations, we calculate the importance score for the linear weight $\mathbf{W}^{l,2}$ in Eq 1 for both target and reference prompts using the methodology described in 4.2 and Eq 2 as

$$\begin{aligned} \mathbf{S}(\mathbf{W}^{l,2}, \mathbf{H}_t^l(\mathcal{P}^*)) &= |\mathbf{W}^{l,2}| \odot (\mathbf{1}^d \cdot \|\mathbf{H}_t^l(\mathcal{P}^*)\|_2) \\ \mathbf{S}(\mathbf{W}^{l,2}, \mathbf{H}_t^l(\mathcal{P})) &= |\mathbf{W}^{l,2}| \odot (\mathbf{1}^d \cdot \|\mathbf{H}_t^l(\mathcal{P})\|_2) \end{aligned} \quad (3)$$

For ease of notation, we denote $\mathbf{S}(\mathbf{W}^{l,2}, \mathbf{H}_t^l(\mathcal{P}^*))$ and $\mathbf{S}(\mathbf{W}^{l,2}, \mathbf{H}_t^l(\mathcal{P}))$ as $\mathbf{S}_t^l(\mathcal{P}^*)$ and $\mathbf{S}_t^l(\mathcal{P})$ respectively in the subsequent sections. Following this, we identify a skilled neuron by comparing its importance score for the target concept prompt with that for the reference prompt.

Isolating concept-generating neurons at time t : Similar to Wanda [Sun et al., 2024], we adopt a *per-output comparison group*, which considers the importance scores among weights in each row of the weight matrix, rather than the matrix as a whole. Specifically, for a given sparsity level $k\%$, we define the top- $k\%$ important weight neurons for generating the target concept in row- i denoted by $\mathbf{W}^{l,2}[i, :]$ as

$$\mathbf{I}_t^l(\mathcal{P}^*)[i, j] = \begin{cases} 1 & \text{if } \mathbf{S}_t^l(\mathcal{P}^*)[i, j] \in \text{top-}k\% \text{ of } \mathbf{S}_t^l(\mathcal{P}^*)[i, :] \\ 0 & \text{otherwise,} \end{cases} \quad (4)$$

where $\mathbf{I}_t^l(\mathcal{P}^*)$ forms a binary mask matrix for the concept prompt set \mathcal{P}^* . As \mathcal{P}^* contains additional undesired target concepts compared with \mathcal{P} , $\mathbf{I}_t^l(\mathcal{P}^*)$ thus consists of the set of important neurons that are responsible for generating both the target and reference concepts. Our next step involves filtering and disentangling these neurons to isolate them to generate the target concept and the reference separately. Continuing with comparison on the Wanda score matrices for both target and reference prompts sets, we now define *skilled* neurons.

Definition 4.1 For a linear layer characterized by $\mathbf{W}^{l,2}$, the weight neuron $\mathbf{W}^{l,2}[i, j]$ is defined as a *skilled neuron* at time step t if $\mathbf{I}_t^l[i, j](\mathcal{P}^*) = 1$ and $\mathbf{S}_t^l(\mathcal{P}^*)[i, j] > \mathbf{S}_t^l(\mathcal{P})[i, j]$.

In essence, if a weight neuron ranks within the top- $k\%$ Wanda scores among other neurons in a row of $\mathbf{W}^{l,2}$ for the target prompts \mathcal{P}^* , it contributes to generating either the undesired target concept or the reference concept. However, if its Wanda score surpasses that of a reference concept, it predominantly influences the target concept.

Subsequently, we form a time-dependent binary mask \mathbf{M}_t^l over weight matrix $\mathbf{W}^{l,2}$ such that

$$\mathbf{M}_t^l[i, j] = \begin{cases} 1 & \text{if weight neuron } \mathbf{W}^{l,2}[i, j] \text{ is skilled} \\ 0 & \text{otherwise,} \end{cases} \quad (5)$$

where \mathbf{M}_t^l is a subset of \mathbf{I}_t^l as only neurons that are highly activated by the target concept are retained.

Removing aggregated skilled neurons over timesteps: While we previously described time-dependent skilled neurons, DiffPrune [Fang et al., 2023] demonstrates that weights can be pruned by aggregating a pruning metric over a selected subset of timesteps based on relative importance scores. However, in our study, we discovered that simply aggregating the binary mask over the first \hat{t} denoising iterations suffices to eliminate a concept while preserving the underlying object. Consequently, we define pruned weight matrix $\hat{\mathbf{W}}^{l,2}$ as

$$\hat{\mathbf{W}}^{l,2} = \mathbf{W}^{l,2} \odot (\neg(\bigvee_{t=T, T-1, \dots, T-\hat{t}} \mathbf{M}_t^l)) \quad (6)$$

where \bigvee and \neg denote the logical OR and NOT operators. All the weights of the pre-trained diffusion model f_θ remain unchanged as only $\mathbf{W}^{l,2}$ is substituted with pruned weights obtained from Equation 6. We then perform experiments with the pruned model to evaluate the effectiveness of concept removal, i.e. subsequently, we only use $\hat{\mathbf{W}}^{l,2}$ for image sampling.

270 5 EXPERIMENTS

271 5.1 EXPERIMENT DETAILS

272 We work with Stable Diffusion v1.5 (SD), which includes 16 FFN layers that serve as candidates for
 273 skilled neuron discovery and pruning. We begin by formulating the calibration sets \mathcal{P}^* and \mathcal{P} that are
 274 used to obtain the matrices $\mathbf{H}_t^l(\mathcal{P}^*)$ and $\mathbf{H}_t^l(\mathcal{P})$ for calculating the score in Equation 3. The list of
 275 prompts and the exact structure of the sentences for different concepts is provided in Table 9 in the
 276 Appendix.

277 **Pruning candidates:** The selection of FFN second layer for pruning was informed by an ablation
 278 study we conducted across various layers within the UNet, aimed at identifying the most effective
 279 pruning targets. Specifically, we analyzed the first layer of the FFNs, along with the query, key, and
 280 value weight matrices in all cross-attention layers. In Appendix A.2, we present the concept erasure
 281 performance for pruning within these layers, as well as an analysis of neuron activation patterns.
 282 The results clearly indicate that the second layer of the FFNs proves to be a more effective pruning
 283 candidate compared to other layers. This observation aligns with findings in the LLM literature,
 284 where these layers have also been identified as prime candidates for skill discovery and pruning
 285 [Zhang et al., 2023e; Suau et al., 2020; Wang et al., 2022]. Finally, to calculate neuron activations,
 286 we run the model for 50 denoising iterations and fix the seed before every forward pass to ensure the
 287 same initializations for both reference and target concept prompts.

288 **Hyper-parameter selection:** As discussed in Section 4.1, we select two key hyperparameters—
 289 sparsity level $k\%$ and \hat{t} —for aggregating skilled neurons across time steps. For each concept,
 290 we vary the sparsity parameter $k\%$ between 0.5% and 5%, choosing the value that achieves the
 291 best trade-off between concept erasure and the retention of unrelated concepts. More details on this
 292 hyperparameter selection process can be found in Section A.3 of the appendix. The optimal sparsity
 293 levels $k\%$ and the corresponding \hat{t} values for each concept are outlined in Table 10 in the appendix.
 294 Interestingly, our experiments reveal that $\hat{t} = 10$ is typically sufficient for erasing concepts while
 295 preserving objects, suggesting that low-level features such as style and objects are formed early in
 296 the denoising process, with fine-grained details being added later.

297 **Baselines:** We identify the following concept editing methods as our closest competitors due to
 298 their lightweight approach: UCE[Gandikota et al., 2023b], Forget-Me-Not (FMN) [Zhang et al.,
 299 2023a], MACE [Lu et al., 2024], Receler [Huang et al., 2023], and SPM [Lyu et al., 2023]. These
 300 works are considered direct competitors as they share a similar emphasis on computational efficiency.
 301 Additionally, we include training-intensive methods such as Concept Ablation (CA) [Kumari et al.,
 302 2023], ESD [Gandikota et al., 2023a], Selective Amnesia (SA)[Heng & Soh, 2023], Scissorhands
 303 (SH)[Wu & Harandi, 2024], and AdvUnlearn [Zhang et al., 2024] in our comparison. However, we
 304 categorize these as indirect competitors, as their reliance on extensive fine-tuning contrasts with
 305 ConceptPrune’s training-free regime. We include a baseline only if their method has been evaluated
 306 for that concept and is reproducible from their source code.¹

307 5.2 ERASING ARTISTIC STYLES

308 We consider five artists — *Van Gogh*, *Claude Monet*, *Pablo Picasso*, *Leonardo Da Vinci*, and *Salvador*
 309 *Dali*. To measure the efficacy of concept removal, we created a dataset of 50 prompts for each artist
 310 using ChatGPT, consisting of the names of their paintings followed by the name of the artist. To
 311 measure the efficacy of concept removal, we report two metrics: the *CLIP Similarity*, which measures
 312 the similarity between the generated image and the prompt, and a stricter *CLIP score* that penalizes a
 313 model when the similarity between the image generated by the concept-editing and the prompt is
 314 greater than the similarity between the image generated by the pre-trained SD and prompt. Lower
 315 values of *CLIP Similarity* and higher values of *CLIP Score* indicate better concept removal. We also
 316 evaluate the fidelity of general purpose image generation by measuring FID and *CLIP Similarity*
 317 on the COCO30k dataset. From the quantitative results presented in Table 2, we demonstrate that
 318 our method outperforms other baselines in artist style removal while effectively retaining unrelated
 319 concepts, as indicated by the low FID score. In Figure 1, we present some qualitative examples
 320 that demonstrate the strong erasing capabilities of ConceptPrune with high-quality realistic output
 321

322 ¹We reproduced CA to remove nudity and object classes from ImageNette but performance was very poor.

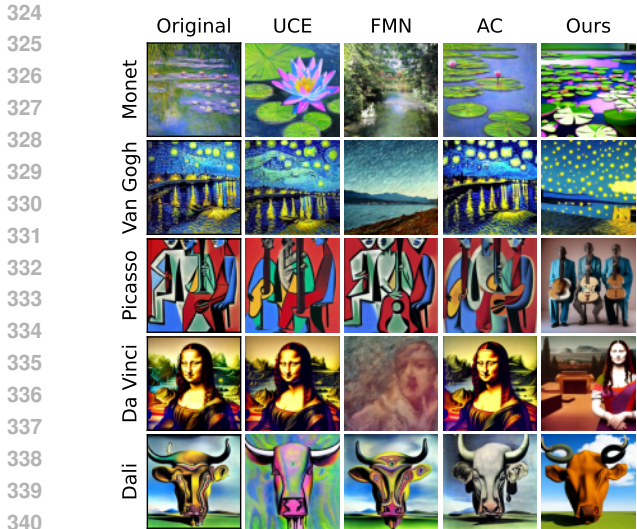


Figure 1: Qualitative results of artist erasure. ConceptPrune demonstrates stronger erasing while generating high-quality, realistic-looking images.

Table 2: Quantitative results of Artist style removal, average over 5 artist styles. CLIP Similarity and CLIP Accuracy measure art style removal. FID and CLIP Similarity on COCO30k measure fidelity for unrelated retained concepts. The full split of the results for different art styles is reported in the appendix in Table 11. Our ConceptPrune can effectively erase artist styles without compromising the model’s performance on unrelated concepts.

Light-weight	Artist erasure			COCO	
	Original SD	Similarity ↓	Score ↑	FID ↓	Similarity ↑
	Original SD	42.1	23.0	14.5	31.3
✗	ESD	34.1	49.2	15.9	30.7
	CA	32.4	65.2	17.5	31.3
	SA	27.1	86.9	14.7	31.3
	AdvUnlearn	27.2	82.0	16.9	29.7
✓	UCE	32.8	44.0	15.7	31.3
	FMN	28.4	82.4	20.9	29.8
	MACE	28.2	85.4	15.1	31.0
	Receler	28.4	82.0	16.7	29.1
✓	Ours	26.9	94.0	16.9	29.9

images. More qualitative results are presented in Section A.4 in the appendix. While we demonstrate strong retention of unrelated concepts in COCO30k in Table 2, Section A.4 in the appendix further provides evidence that using ConceptPrune to erase an artist’s style results in minimal degradation when generating other similar artist styles.

5.3 ERASING EXPLICIT CONTENT

We quantitatively evaluate our proposed method for moderating Not-Safe-for-Work (NSFW) concepts like nudity by comparing it against the concept-erasing baselines ESD, UCE, and FMN. In addition, we also compare with variants of Stable Diffusion, such as Safe Latent Diffusion (SLD) [Schramowski et al., 2023] and Stable Diffusion 2.0 [Rombach et al., 2022a], which have been fine-tuned on a filtered subset of LAION without explicit images. We use the Inappropriate Prompts Dataset (I2P) [Schramowski et al., 2023], which consists of 4703 prompts featuring various inappropriate concepts. Nudity detectors [Bedapudi, 2022] indicate that, out of these 4703 prompts, the pre-trained Stable Diffusion model generates nudity for 796 prompts. In Figure 2, we report the percentage reduction in the number of generated images with nudity compared to the pre-trained Stable Diffusion model. ConceptPrune generates nudity in merely 47 prompts within 4703 prompts in the I2P dataset, implying a 94.1% decrease compared to 88% in ESD and 85.6% in UCE, demonstrating a significant improvement over other baselines in content moderation. We present more qualitative results on the I2P dataset in Figure 13 in the appendix.

5.4 ERASING OBJECTS

Single-object erasing: We showcase the effectiveness of our method in removing objects from the learned concepts of diffusion models. We conducted experiments targeting ImageNette classes [Howard & Gugger, 2020], a subset of ImageNet [Deng et al., 2009] comprising 10 classes. Similar to UCE and ESD, we generated 500 images per class and evaluated the top-1 classification accuracy using a pre-trained ResNet-50 [He et al., 2015]. Table 4 shows that ConceptPrune has superior erasure performance on average while effectively minimizing interference on non-targeted classes. ² More results of object erasure are provided in Figure 14 in the appendix.

²We copied the numbers from their original papers based on SD 1.4 and therefore, we repeated our experiments with SD 1.4 for consistency.

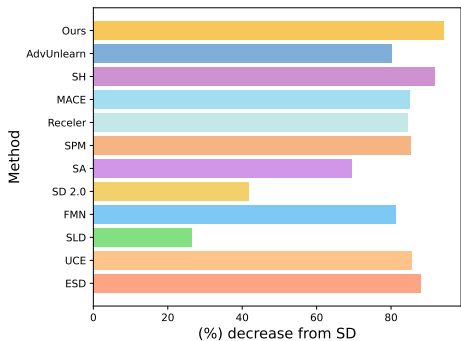


Figure 2: Explicit Content Erasure. The percentage reduction in nudity content from I2P prompts, compared to the original SD model ConceptPrune (SD1.5) decreases the number of explicit images by 94.1%, outperforming competitors as well as SD2.0.

Table 3: ConceptPrune demonstrates robustness to adversarial attacks. Unlearning methods evaluated against three adversarial attacks. Black-box (Ring-A-Bell[Tsai et al., 2024], and MMA[Yang et al., 2023]) performance is quantified by percentage reduction in nude samples compared to SD. White-box UnlearnDiffAtk [Zhang et al., 2023d] performance measures the attack success rate (ASR).

Light-weight	Method	Ring-A-Bell \uparrow	MMA \uparrow	UnlearnDiffAtk \downarrow
-	SLD	2.8	25.5	82.4
	SDv2	1.8	26.8	73.8
x	ESD	52.8	87.3	76.1
	SA	84.3	94.3	11.3
	SH	86.1	94.3	22.3
	AdvUnlearn	85.8	93.7	21.1
	UCE	67.6	63.3	93.2
✓	Receler	67.9	65.7	92.1
	MACE	56.4	57.9	89.3
	FMN	5.6	53.6	97.9
	SPM	34.5	78.4	91.6
✓	Ours	85.2	95.6	64.8

Table 4: Concept Erasure: Top-1 classification accuracy of erased and preserved class samples, using a pre-trained ResNet-50. Our ConceptPrune effectively erases objects from pre-trained models without impacting the accuracy for other object classes.

Classes	Accuracy of Erased Classes \downarrow				Accuracy of Preserved Classes \uparrow			
	ESD	UCE	FMN	ConceptPrune	ESD	UCE	FMN	ConceptPrune
Church	54.2	8.4	2.0	6.0	80.2	77.5	57.8	82.8
English Springer	6.2	0.2	1.9	0.0	62.6	78.9	73.5	80.1
Golf ball	5.8	0.8	13.7	0.0	65.6	79.0	82.8	87.8
Gas Pump	8.6	0.0	7.9	0.0	66.5	80.7	79.0	83.0
Tench	9.6	0.0	5.7	0.0	66.6	79.3	78.4	85.0
Parachute	23.8	1.4	8.3	7.0	65.4	77.4	98.2	80.6
Cassette Player	0.6	0.0	1.0	1.0	64.5	90.3	68.7	94.3
Chain Saw	6.0	0.0	0.1	0.0	71.6	80.2	78.4	91.5
French Horn	0.4	3.0	0.0	3.0	77.0	80.1	78.3	88.2
Garbage Truck	10.4	14.8	0.1	0.0	51.5	78.7	74.9	85.8
Average	12.5	2.7	4.1	1.7	66.9	80.2	77.5	85.9

Multi-object erasing: In addition to single-object erasing, we also evaluate ConceptPrune on removing multiple objects from the model simultaneously. Although our pruning strategy generates a pruning mask for concepts individually, it provides a straightforward baseline for multi-object erasing by taking the union of skilled neurons across different concepts. We direct the reader to Appendix A.5 for more details. We compare our method with UCE and report the accuracy on erased classes along with FID and CLIP similarity on COCO30k. Table 6 shows that ConceptPrune demonstrates comparable erasing performance while excelling at retaining unrelated concepts.

5.5 ADVERSARIAL DEFENSE ON CONCEPT ERASURE ATTACKS

White-box attacks: Recent research has recognized the limitations of the concept editing baselines considered in this paper, namely UCE, ESD, FMN, and CA. Model-based adversarial attacks like UnlearnDiffAtk [Zhang et al., 2023d] and Concept Inversion (CI) [Pham et al., 2023] have demonstrated that subtle perturbations to text prompts can circumvent the unlearning mechanisms, compelling concept-editing baselines to generate harmful images with undesired concepts once again. Furthermore, these studies show a near-perfect Attack Success Rate (ASR) for FMN and UCE which jeopardizes the safety and effectiveness of these baselines in real-world settings.

Table 5: ConceptPrune is substantially more robust to adversarial attacks aimed at eliciting erased concepts. **(Top):** Attack Success Ratio (ASR %, ↓) of UnlearnDiffAtk [Zhang et al., 2023d] adversarial prompts for Van Gogh’s painting style and 4 classes of the Imagenette dataset.

Light-weight	Artist Style		Object erasing				
	Vincent Van Gogh Top-1 ASR	Vincent Van Gogh Top-3 ASR	Parachute ASR	Tench ASR	Garbage Truck ASR	Church ASR	
×	ESD	32.0	76.0	54.0	36.0	24.0	60.0
	CA	77.0	92.0	–	–	–	–
	SH	–	–	24.0	8.0	2.0	6.0
	SalUn	–	–	74.0	14.0	42.0	62.0
	AdvUnlearn	2.0	24.6	14.0	4.0	8.0	6.0
✓	UCE	94.0	100.0	43.0	22.0	38.0	68.0
	FMN	56.0	90.0	100.0	100.0	98.0	96.0
	SPM	–	–	96.0	90.0	82.0	94.0
✓	ConceptPrune (Ours)	2.4	24.4	34.0	16.1	0.0	21.7

Table 6: Quantitative results for multi-object erasure. We report Accuracy on erased classes and FID on COCO30k, CLIP similarity on COCO30k, and ASR of UnlearnDiffAtk. ConceptPrune is comparable to UCE at erasing multiple objects and outperforms UCE in retaining image generation capabilities along with being significantly robust to white-box adversaries.

	COCO FID	CLIP score	Accuracy on erased classes	ASR
UCE	17.7	31.0	4%	22%
ConceptPrune	17.5	29.9	7%	6%

We evaluate ConceptPrune under these recently introduced white-box attacks - UnlearnDiffAtk [Zhang et al., 2023d] and Concept Inversion (CI) [Pham et al., 2023]. For UnlearnDiffAtk, we evaluate for Van Gogh style, ImageNette objects, and nudity. We compare ConceptPrune with baselines UCE, ESD, and FMN across all concepts, and for nudity, we include comparisons with presumably safe models such as Safe Latent Diffusion (SLD) and SDv2. Following [Zhang et al., 2023d], we report the top-1 and top-3 ASR for Van Gogh style, which indicates whether the generated image is classified as the top-1 prediction or within the top-3 predictions for Van Gogh’s painting style when evaluated by the post-generation image classifier. For object erasure and NSFW attacks, we report ASR based on a pre-trained ResNet50 model and NudeNet detector respectively [Bedapudi, 2022]. Table 5 (top) illustrates that for artist style and object erasure, ConceptPrune renders the UnlearnDiffAtk unsuccessful, achieving a 0% ASR in two instances, in contrast to the perfect success rates seen for baselines like UCE and FMN. Table 3 shows that UCE, ESD, and FMN fail to defend against the NSFW attack, ConceptPrune demonstrates an ASR of 64.8%, significantly lower than that the models that are trained for safety (SDv2 and SLD).

Following the evaluation protocol of Concept Inversion (CI), we generated 500 images per class and evaluated the top-1 classification accuracy. Similar to CI, we also compare the performance of ConceptPrune against negative prompting (Neg-Prompt) [Yuanhao et al., 2024] and Safe Latent Diffusion (SLD-Med) [Schramowski et al., 2023]. In Table 13 in the Appendix, we observe that the accuracy of 3 out of 4 erased classes is notably lower compared to other baselines. This demonstrates that ConceptPrune offers significantly greater adversarial robustness against various white-box attack variants. We present more qualitative analysis in Figure 11 in the appendix.

Black-box attacks: To prevent the generation of NSFW imagery, SD models incorporate preventive measures such as prompt filters and post-synthesis safety checks by default. In a black-box setting such as a web service, these defenses are considered impossible to override. Therefore, we also evaluate black-box robustness. Recent research MMA-Diffusion [Yang et al., 2023] released a set of 1000 adversarial prompts for SDv1.5 that circumvent safety filters on the text and image level. In addition, Ring-A-Bell [Tsai et al., 2024] directly challenges our competitors ESD, UCE, and FMN and attacks their erasing strength with their set of adversarial prompts. Inspired by these works, we evaluate ConceptPrune along with competitors on adversarial prompts released by [Yang et al., 2023; Tsai et al., 2024] and report the percentage reduction in number of images for which nudity is

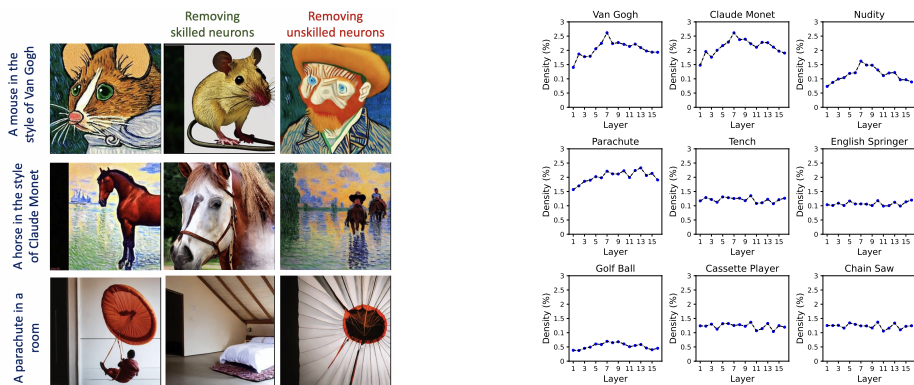


Figure 3: *Left*: ConceptPrune effectively disentangles skilled neurons responsible for specific concepts from general object-generating neurons. E.g., removing "Van Gogh" skilled neurons erases the "Van Gogh" style while removing unskilled neurons eliminates the object. *Right*: Skilled neurons are localized to a very compact subspace, between 1% to 3% of FFN parameters.

generated as compared to pre-trained SD. Results in Table 3 show that ConceptPrune offers a stark increase in adversarial robustness with a 95.6% decrease in the generation of nudity under MMA. This underscores its potential as a reliable and safe choice over our competitors. We present more qualitative analysis in Figure 11 in the appendix.

5.6 FURTHER ANALYSIS

Analysing the density of skilled neurons: We evaluate the *density* of skilled neurons, defined as the percentage of non-zero elements in the pruning mask in Equation 5. Our analysis in Figure 3 (right) reveals that concept-generating neurons span less than 3% of the FFN weights matrix considered for pruning. This suggests that concept generation can be attributed to a very tiny subspace, potentially constituting less than 0.12% of the total model parameters in diffusion models.

Are concept-generating skilled neurons disentangled from object-generating neurons? In Section 5, we demonstrated that ConceptPrune exhibits strong concept erasure skills for a diverse range of concepts by discovering and pruning a compact subspace of skilled neurons. Conversely, removing unskilled neurons, i.e. neurons that satisfy the opposite of the second condition in Definition 4.1 and follow $\mathbf{S}_t^l(\mathcal{P}^*)[i, j] < \mathbf{S}_t^l(\mathcal{P})[i, j]$ instead are hypothesised to distort the reference concept while retaining the target concept. Figure 3 (left) offers qualitative examples that confirm our hypothesis, illustrating our ability to isolate a distinct set of neurons solely responsible for generating concepts, demonstrating their disentanglement from neurons responsible for generating general utilities. We present an interesting study on gender-specific neurons in diffusion models in Section A.7.

Can ConceptPrune generalize to other architectures? We demonstrate that ConceptPrune can be seamlessly applied to Stable Diffusion v2.0 and SD-XL. We erased the artist styles listed in Table 2 and compared the results with UCE on SD-v2.0 and SD-XL. As shown in Table 14 in the appendix, ConceptPrune not only generalizes well to different architectures but also delivers superior erasure performance across models.

6 CONCLUSIONS

This paper revisited the important challenge of concept editing in pre-trained diffusion models from the perspective of skilled neuron identification and pruning. We showed that concepts related to object categories, art styles, gender, and nudity can be identified and pruned – leading to effective erasure while maintaining overall generation quality. Our ConceptPrune approach is fast, training-free, and permanent – exhibiting strong robustness to adversarial attacks that break prior concept erasure methods. Without relying on token-rewriting, pruned models could be distributed without the risk of adversaries simply removing rewriting safeguards. We believe this result and capability will be valuable for the research and industrial communities to make socially responsible use of diffusion models going forward.

REFERENCES

- 540
541
542 Omer Antverg and Yonatan Belinkov. On the pitfalls of analyzing individual neurons in language
543 models. *arXiv*, 2022.
- 544 Praneeth Bedapudi. Nudenet: Neural nets for nudity detection and censoring. 2022.
- 545
546 Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Guttag. What is the state of
547 neural network pruning? *arXiv*, 2020.
- 548 Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim.
549 Coyo-700m: Image-text pair dataset. 2022.
- 550
551 Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing
552 web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021.
- 553 Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James Glass. What
554 is one grain of sand in the desert? analyzing individual neurons in deep nlp models. *arXiv*, 2018.
- 555
556 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: a large-scale
557 hierarchical image database. pp. 248–255, 06 2009. doi: 10.1109/CVPR.2009.5206848.
- 558
559 Nadir Durrani, Hassan Sajjad, Fahim Dalvi, and Yonatan Belinkov. Analyzing individual neurons
560 in pre-trained language models. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.),
561 *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*
562 (*EMNLP*), pp. 4865–4880, Online, November 2020. Association for Computational Linguistics.
563 doi: 10.18653/v1/2020.emnlp-main.395. URL [https://aclanthology.org/2020.
564 emnlp-main.395](https://aclanthology.org/2020.emnlp-main.395).
- 565
566 Nadir Durrani, Fahim Dalvi, and Hassan Sajjad. Discovering salient neurons in deep nlp models.
567 *Journal of Machine Learning Research*, 24(362):1–40, 2023.
- 568
569 Gemini Team et al. Gemini: A family of highly capable multimodal models. *arXiv*, 2024.
- 570
571 Sarah Andersen. et al v. Stability AI Ltd. et al. Case no.3:2023cv00201. us district court for the
572 northern district of california., 2023.
- 573
574 Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. Salun: Em-
575 powering machine unlearning via gradient-based weight saliency in both image classification and
576 generation. *arXiv preprint arXiv:2310.12508*, 2023.
- 577
578 Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. *arXiv*, 2023.
- 579
580 Camera Forensics. The dark reality of stable diffusion. 2024.
- 581
582 Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural
583 networks. *arXiv*, 2019.
- 584
585 Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in
586 one-shot. *arXiv*, 2023.
- 587
588 Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training
589 quantization for generative pre-trained transformers. *arXiv*, 2023.
- 590
591 Rohit Gandikota, Joanna Materzyńska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts
592 from diffusion models. In *Proceedings of the 2023 IEEE International Conference on Computer
593 Vision*, 2023a.
- 594
595 Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified
596 concept editing in diffusion models. *arXiv preprint arXiv:2308.14761*, 2023b.
- 597
598 Kristian Georgiev, Joshua Vendrow, Hadi Salman, Sung Min Park, and Aleksander Madry. The
599 journey, not the destination: How data guides diffusion models. In *Arxiv preprint arXiv:2312.06205*,
600 2023.

- 594 Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for
595 efficient neural networks. *arXiv*, 2015.
- 596 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
597 recognition. *arXiv*, 2015.
- 598
- 599 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv*, 2023.
- 600
- 601 Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep
602 generative models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
603 URL <https://openreview.net/forum?id=BC1IJdsuYB>.
- 604 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint*
605 *arxiv:2006.11239*, 2020.
- 606
- 607 Jeremy Howard and Sylvain Gugger. Fastai: A layered api for deep learning. *Information*, 11(2):108,
608 February 2020. ISSN 2078-2489. doi: 10.3390/info11020108. URL [http://dx.doi.org/](http://dx.doi.org/10.3390/info11020108)
609 [10.3390/info11020108](http://dx.doi.org/10.3390/info11020108).
- 610 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
611 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International*
612 *Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=nZeVKeeFYf9)
613 [id=nZeVKeeFYf9](https://openreview.net/forum?id=nZeVKeeFYf9).
- 614 Chi-Pin Huang, Kai-Po Chang, Chung-Ting Tsai, Yung-Hsuan Lai, Fu-En Yang, and Yu-Chiang Frank
615 Wang. Receler: Reliable concept erasing of text-to-image diffusion models via lightweight erasers.
616 *arXiv preprint arXiv:2311.17717*, 2023.
- 617
- 618 Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu.
619 Ablating concepts in text-to-image diffusion models. In *International Conference on Computer*
620 *Vision (ICCV)*, 2023.
- 621 Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In D. Touretzky
622 (ed.), *Advances in Neural Information Processing Systems*, volume 2. Morgan-Kaufmann,
623 1989. URL [https://proceedings.neurips.cc/paper_files/paper/1989/](https://proceedings.neurips.cc/paper_files/paper/1989/file/6c9882bbac1c7093bd25041881277658-Paper.pdf)
624 [file/6c9882bbac1c7093bd25041881277658-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1989/file/6c9882bbac1c7093bd25041881277658-Paper.pdf).
- 625 Namhoon Lee, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Snip: Single-shot network pruning
626 based on connection sensitivity. *arXiv*, 2019.
- 627
- 628 Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for
629 efficient convnets. *arXiv*, 2017.
- 630 Zhili Liu, Kai Chen, Yifan Zhang, Jianhua Han, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan
631 Yeung, and James Kwok. Implicit concept removal of diffusion models. *arXiv*, 2024.
- 632
- 633 Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value
634 of network pruning. In *International Conference on Learning Representations*, 2019. URL
635 <https://openreview.net/forum?id=rJlnB3C5Ym>.
- 636
- 637 Shilin Lu, Zilan Wang, Leyang Li, Yanzhu Liu, and Adams Wai-Kin Kong. Mace: Mass concept
638 erasure in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
Pattern Recognition, pp. 6430–6440, 2024.
- 639
- 640 Sasha Luccioni, Christopher Akiki, Margaret Mitchell, and Yacine Jernite. Stable bias: Evaluating
641 societal representations in diffusion models. In *Thirty-seventh Conference on Neural Information*
642 *Processing Systems Datasets and Benchmarks Track*, 2023. URL [https://openreview.](https://openreview.net/forum?id=qVXYU3F017)
[net/forum?id=qVXYU3F017](https://openreview.net/forum?id=qVXYU3F017).
- 643
- 644 Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models:
645 Synthesizing high-resolution images with few-step inference. *arXiv*, 2023.
- 646
- 647 Mengyao Lyu, Yuhong Yang, Haiwen Hong, Hui Chen, Xuan Jin, Yuan He, Hui Xue, Jungong Han,
and Guiguang Ding. One-dimensional adapter to rule them all: Concepts, diffusion models and
erasing applications, 2023.

- 648 Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting
649 them. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
650 5188–5196, 2015.
- 651 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
652 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
653 text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 654 Hadas Orgad, Bahjat Kawar, and Yonatan Belinkov. Editing implicit assumptions in text-to-image
655 diffusion models. *arXiv:2303.08084*, 2023.
- 656 Minh Pham, Kelly O. Marshall, and Chinmay Hegde. Circumventing concept erasure methods for
657 text-to-image generative models. *arXiv*, 2023.
- 658 Minh Pham, Kelly O. Marshall, Niv Cohen, Govind Mittal, and Chinmay Hegde. Circumventing
659 concept erasure methods for text-to-image generative models. In *The Twelfth International
660 Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?
661 id=ag3o2T51Ht](https://openreview.net/forum?id=ag3o2T51Ht).
- 662 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
663 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
664 synthesis. *arXiv*, 2023.
- 665 MIT Technology Review. Text-to-image ai models can be tricked into generating disturbing images.
666 2023.
- 667 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
668 resolution image synthesis with latent diffusion models. *arXiv*, 2021.
- 669 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
670 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-
671 ence on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, June 2022a.
- 672 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
673 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-
674 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022b.
- 675 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical
676 image segmentation. *arXiv*, 2015.
- 677 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.
678 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. 2022.
- 679 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed
680 Kamyar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans,
681 Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion
682 models with deep language understanding. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,
683 and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL
684 <https://openreview.net/forum?id=08Yk-n512A1>.
- 685 Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion:
686 Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE Conference
687 on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- 688 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi
689 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,
690 Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia
691 Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models.
692 In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks
693 Track*, 2022. URL <https://openreview.net/forum?id=M3Y74vmsMcY>.
- 694 Noam Shazeer. Glu variants improve transformer. *arXiv*, 2020.

- 702 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Inter-*
703 *national Conference on Learning Representations*, 2021. URL [https://openreview.net/](https://openreview.net/forum?id=StlgiaRCHLP)
704 [forum?id=StlgiaRCHLP](https://openreview.net/forum?id=StlgiaRCHLP).
- 705 Xavier Suau, Luca Zappella, and Nicholas Apostoloff. Finding experts in transformer models. *arXiv*
706 *preprint arXiv:2005.07647*, 2020.
- 707 Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. A simple and effective pruning approach
708 for large language models. *arXiv*, 2024.
- 709 Yu-Lin Tsai, Chia-Yi Hsu, Chulin Xie, Chih-Hsun Lin, Jia-You Chen, Bo Li, Pin-Yu Chen, Chia-Mu
710 Yu, and Chun-Ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion
711 models? *arXiv*, 2024.
- 712 Xiazhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. Finding skill
713 neurons in pre-trained transformer-based language models. *arXiv preprint arXiv:2211.07349*,
714 2022.
- 715 Boyi Wei, Kaixuan Huang, Yangsibo Huang, Tinghao Xie, Xiangyu Qi, Mengzhou Xia, Prateek
716 Mittal, Mengdi Wang, and Peter Henderson. Assessing the brittleness of safety alignment via
717 pruning and low-rank modifications. *arXiv preprint arXiv:2402.05162*, 2024.
- 718 Jing Wu and Mehrtash Harandi. Scissorhands: Scrub data influence via connection sensitivity in
719 networks. *arXiv preprint arXiv:2401.06187*, 2024.
- 720 Jing Wu, Trung Le, Munawar Hayat, and Mehrtash Harandi. Erasediff: Erasing data influence in
721 diffusion models. 2024. URL <https://openreview.net/forum?id=eVpjeCNsR6>.
- 722 Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Nan Xu, and Qiang Xu. Mma-diffusion: Multimodal
723 attack on diffusion models. *arXiv preprint arXiv:2311.17516*, 2023.
- 724 Ban Yuanhao, Wang Ruochen, Zhou Tianyi, Cheng Minhao, Gong Boqing, and Hsieh Cho-Jui.
725 Understanding the impact of negative prompts: When and how do they take effect? In *arXiv*
726 *preprint*, 2024.
- 727 Eric Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not: Learning
728 to forget in text-to-image diffusion models. *arXiv preprint arXiv:2211.08332*, 2023a.
- 729 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
730 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*
731 *(ICCV)*, pp. 3836–3847, October 2023b.
- 732 Tianyi Zhang, Zheng Wang, Jing Huang, Mohiuddin Muhammad Tasnim, and Wei Shi. A sur-
733 vey of diffusion based image generation models: Issues and their solutions. *arXiv preprint*
734 *arXiv:2308.13142*, 2023c.
- 735 Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and
736 Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate
737 unsafe images... for now. *arXiv preprint arXiv:2310.11868*, 2023d.
- 738 Yimeng Zhang, Xin Chen, Jinghan Jia, Yihua Zhang, Chongyu Fan, Jiancheng Liu, Mingyi Hong,
739 Ke Ding, and Sijia Liu. Defensive unlearning with adversarial training for robust concept erasure
740 in diffusion models. *arXiv preprint arXiv:2405.15234*, 2024.
- 741 Zhengyan Zhang, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. Moefication:
742 Transformer feed-forward layers are mixtures of experts. *arXiv*, 2022.
- 743 Zhengyan Zhang, Zhiyuan Zeng, Yankai Lin, Chaojun Xiao, Xiazhi Wang, Xu Han, Zhiyuan Liu,
744 Ruobing Xie, Maosong Sun, and Jie Zhou. Emergent modularity in pre-trained transformers. *arXiv*
745 *preprint arXiv:2305.18390*, 2023e.
- 746 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in
747 coreference resolution: Evaluation and debiasing methods. *arXiv*, 2018.
- 748 Mengnan Zhao, Lihe Zhang, Tianhang Zheng, Yuqiu Kong, and Baocai Yin. Separable multi-concept
749 erasure from diffusion models. *arXiv*, 2024.

A APPENDIX

A.1 LIMITATIONS

While erasing specific objects, such as the "English Springer," we noticed that a few related dog breeds were also inadvertently removed. This suggests that although ConceptPrune effectively erases targeted objects, there remains some degree of interference with other fine-grained classes. Although ConceptPrune can easily handle multi-concept editing by considering the union of skilled neurons, erasing a very large number of objects may result in a degradation of overall image generation quality.

A.2 SELECTION OF PRUNING CANDIDATES

In this section, we conduct an ablation study on various candidate layers within the UNet to determine the most effective pruning targets. Specifically, we examined the first layer of the FFN (**FFN-1**), second layer of the FFN (**FFN-2**), the Key weight matrix in the Cross Attention layer (**CA-Key**), the Value weight matrix in the Cross Attention layer (**CA-Value**), and the second layer of the FFNs in the text encoder (**CLIP**). **CA-Key** and **CA-Value** were considered because these weight matrices operate on text tokens, while the noised latent tokens are used as queries. We then apply ConceptPrune for pruning different parameters within these layers and report the concept erasure performance in Tables 7 and 8. Firstly, we visually observed that pruning **CA-Key** degrades image quality by distorting objects and textures. Therefore, we have decided not to report the erasure performance associated with pruning **CA-Key**. From Tables 7 and 8, we empirically observed that **FFN-2** is the best choice for pruning.

Additionally, we analyzed neuron activation patterns of different layers to understand which layers consist of neurons that are indicative of the presence of a particular concept. For a given layer, we calculate the norm of activations for input neurons over reference and target prompts, averaging over denoising time steps. The top 1% of neurons in UNet are then identified and their distribution is plotted in 4 (b, c, d). We observed a significant difference in distributions' means in the 2nd FFN layer, indicating distinct activations for reference and target prompts. This distinction is absent in other layers. From these results, it is evident that FFN-2 is a better and sensible pruning candidate than others.

Table 7: Accuracy of erased classes (\downarrow) and preserved classes (\uparrow) for object erasure across different pruning candidates. FFN-2 is a better pruning target.

Pruning candidate	FFN-2(in the paper)		FFN-1		CA-Value		CLIP	
	Erased	Preserved	Erased	Preserved	Erased	Preserved	Erased	Preserved
Parachute	6.9	72.8	21.0	62.2	32.0	69.2	38.0	47.8
English springer	0.0	93.7	46.2	90.0	32.8	89.2	1.0	42.3
French horn	1.9	74.5	17.0	74.8	31.4	79.2	18.0	72.4
Tench	0.0	90.1	47.0	87.1	21.2	73.4	39.0	89.2

Table 8: Erasure performance for artist style removal (first 5 rows, CLIP similarity between the generated image and prompt (\downarrow)) and Nudity (last row, % nudity reduction (\uparrow)) across different pruning candidates. The sparsity level used for pruning is 2%. FFN-2 is a better pruning candidate.

Pruning candidate	Van Gogh	Monet	Leonardo Da Vinci	Pablo Picasso	Salvador Dali	Nudity
FFN-2 (in the paper) (2%)	29.2	23.6	26.5	25.3	29.8	94.1
FFN-1 (2%)	32.7	30.6	29.0	26.5	30.7	67.8
CA-Value (2%)	32.7	30.3	28.6	27.7	27.7	46.2
CLIP (2%)	33.2	32.6	29.4	28.7	31.7	9.1

A.3 DETAILS ON PROMPTS AND HYPER-PARAMETERS

Selecting optimal sparsity ratio - To understand the effect of sparsity level (k%), we vary it from 0.5% to 5% and plot erasure vs. retention performance. Erasure performance is measured by the CLIP similarity between the generated image and the input prompt, with lower values indicating better erasure. Retention is evaluated using a subset of COCO dataset prompts, measuring CLIP similarity between the generated image and the input prompt. From Figure 4(a), we observed that

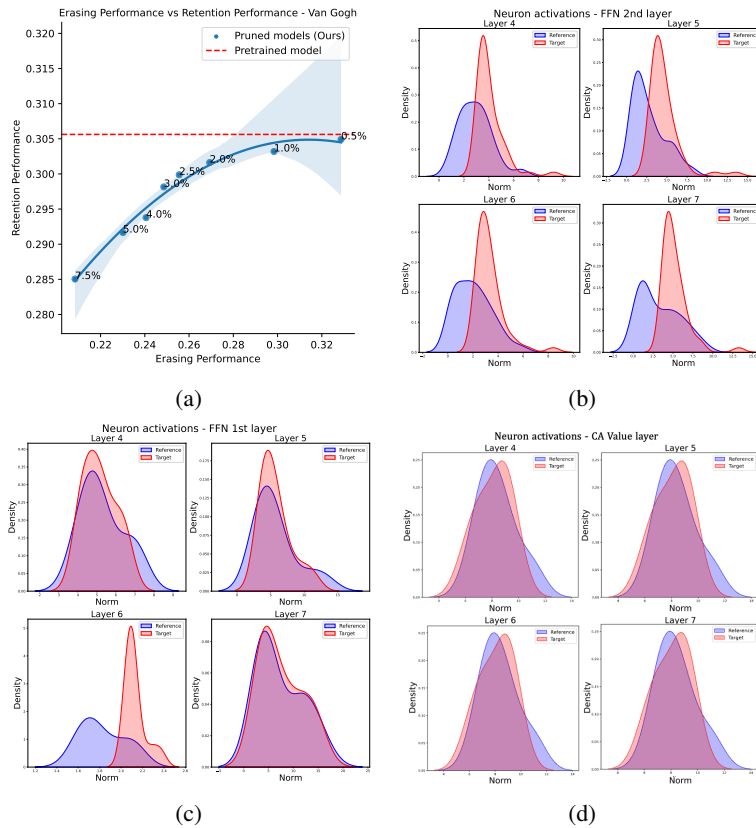


Figure 4: (a): Erasing vs Retention Performance with varying sparsity thresholds. Concept - Van Gogh. (b, c, d): Density of neuron activations for reference and target prompts in second layer of FFNs (pruned in paper), first layer of FFNs and Value layer in cross-attentions respectively. FFN-2 has the most distinct activation distribution.

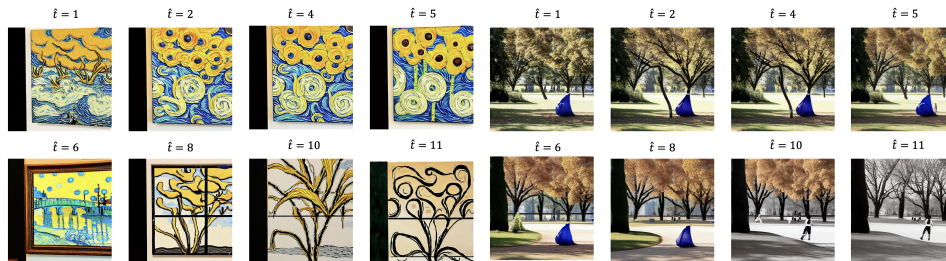


Figure 5: We present qualitative results by varying \hat{t} from 1 to 15 and visualizing the images after concept erasure. Extending beyond 10 timesteps results in a noticeable degradation of image quality.

a sparsity level of $k = 2.5\%$ or $k = 2\%$ offers a good balance of improved erasure with a minimal retention loss (main experiments used 2%).

Selecting optimal \hat{t} - As noted in Section 4, our work draws inspiration from the study in DiffPrune [Fang et al., 2023], which utilizes Taylor expansion at pruned timesteps to estimate weight importance. Their findings reveal that earlier timesteps focus on local features like edges and colors, while later timesteps shift attention to broader content, such as objects and shapes. Similar to [Fang et al., 2023], [Georgiev et al., 2023] also shows that properties such as background color, object shape, etc generated in the earlier timesteps are carried forward to later in the denoising trajectory. Since our work primarily addresses local properties such as style, color, and object shape, we focus on

removing concept-generating neurons in the earlier timesteps, which contribute the most to salient properties in the image.

To determine \hat{t} in Equation 5, we performed a straightforward grid search over timesteps $t = 1$ to $t = 15$, examining erased images to identify the timestep at which the concept was effectively removed from the majority of images. We present qualitative results in Figure 5 by varying \hat{t} from 1 to 15 and visualizing the images after concept erasure. The ablation study is conducted on two distinct concepts: the *Van Gogh* artist style and the object *parachute*. The prompts used for this analysis are "A painting in the style of Van Gogh" for the artist's style and "A parachute in a park" for the object. By varying \hat{t} , we visualize the resulting images to evaluate the effectiveness of concept removal at different timesteps. Our findings demonstrate that focusing exclusively on neurons too early in the denoising trajectory fails to capture all the neurons responsible for generating the target concept. On the other hand, extending beyond 10 timesteps results in a noticeable degradation of image content and quality, striking a delicate balance between effective concept removal and preserving the overall integrity of the image. Therefore, $\hat{t} = 10$ is an optimal point for concept erasure and good retention.

Table 9: Details on calibration set formulation for different concepts. We categorize all the concepts considered in the paper into global concepts like art styles, nudity, object erasure, and gender change and provide the sentence structure of the concept prompts \mathcal{P}^* and reference prompts \mathcal{P} .

Global Concept	Concept (<concept>)	Concept prompt	Reference prompt	List
Art Styles	Van Gogh Monet	a <object> in the style of <concept>	a <object>	<object> ∈ [cat, dog, mouse, bear, car, chair, bag, ball, door, rabbit, flower, bowl, bat, horse, bird, wall, bottle, building, chicken, duck]
	Leonardo Da Vinci Salvador Dali Pablo Picasso			
Nudity	naked	a photo of a naked <person>	a photo of a <person>	<person> ∈ A list of person related words ³
Object Erasure	parachute, gas pump golf ball, cassette player	a <concept> in a <scene>	a <scene>	a <scene> ∈ [road, garden, beach room, park, table bag, tree, forest street, shelter, chair]
	english springer, tench chain saw, french horn			
Object Erasure	church, garbage truck	a <concept> near a <place>	a <place>	<place> ∈ road, park, beach, street house, tree, forest, statue, car]
Gender change	Male to Female	a photo of a <male>	a photo of a <female>	<male> ∈ [man, boy, person, guy father, son, husband, uncle]
	Female to Male	a photo of a <female>	a photo of a <male>	<female> ∈ [woman, girl, female, lady mother, daughter, wife, aunt]

A.4 ARTIST STYLE ERASURE

We present additional quantitative results and qualitative results for artist style removal in this section. Please see Figure 6, 7, 8, 9, and 10 and Table 11.

Cross-artist erasure: Ideally, erasing an artist’s style should not impact the generation of other artist styles. However, concept erasure baselines like CA [Kumari et al., 2023] and UCE [Gandikota et al., 2023b] have reported slight degradation in generating paintings of other artists when a similar style is removed. For instance, [Kumari et al., 2023] demonstrates that removing 'Van Gogh' style results in the removal of the 'Claude Monet' style. To assess this quantitatively, we used CA, UCE, and ConceptPrune to erase the 'Van Gogh' style and evaluated the performance over the remaining four artist styles in Table 11. We measure the CLIP similarity between the generated image and the input prompt, where a higher CLIP similarity indicates better preservation of the artist’s style. Table 12 demonstrates that while ConceptPrune performs comparably to other baselines in preserving related artist styles, it outperforms them in maintaining the model’s overall image generation capabilities (Table 2).

A.5 MULTI-OBJECT ERASING

We outline our approach to multi-object erasing, where we take the union of skilled neurons across all targeted objects and prune them collectively. Let the binary mask representing skilled neurons for a concept c in Equation 6 be $\mathbf{M}_c^{t,l}$. For erasing a set of multiple concepts $\mathbb{C} = \{c_1, c_2, \dots, c_m\}$, we take the union of skilled neurons for each time step and concept $\bigvee_{c \in \mathbb{C}} \mathbf{M}_c^{t,l}$, and formulate the pruned matrix $\hat{\mathbf{W}}_l^2$ as $\mathbf{W}_l^2 \odot \left(\neg(\bigvee_{t=T, T-1, \dots, T-\hat{t}} \bigvee_{c \in \mathbb{C}} \mathbf{M}_c^{t,l}) \right)$, where \bigvee and \neg denote the logical OR and NOT operators.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941



Figure 6: Qualitative results for erasing artist - *Van Gogh*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist’s style.

942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

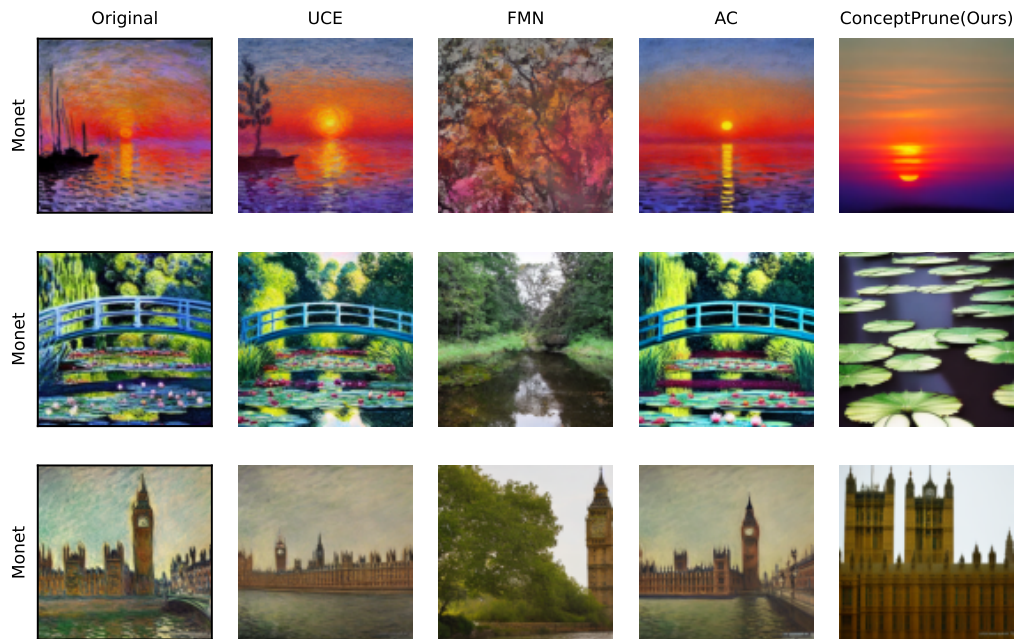


Figure 7: Qualitative results for erasing artist - *Monet*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist’s style.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995

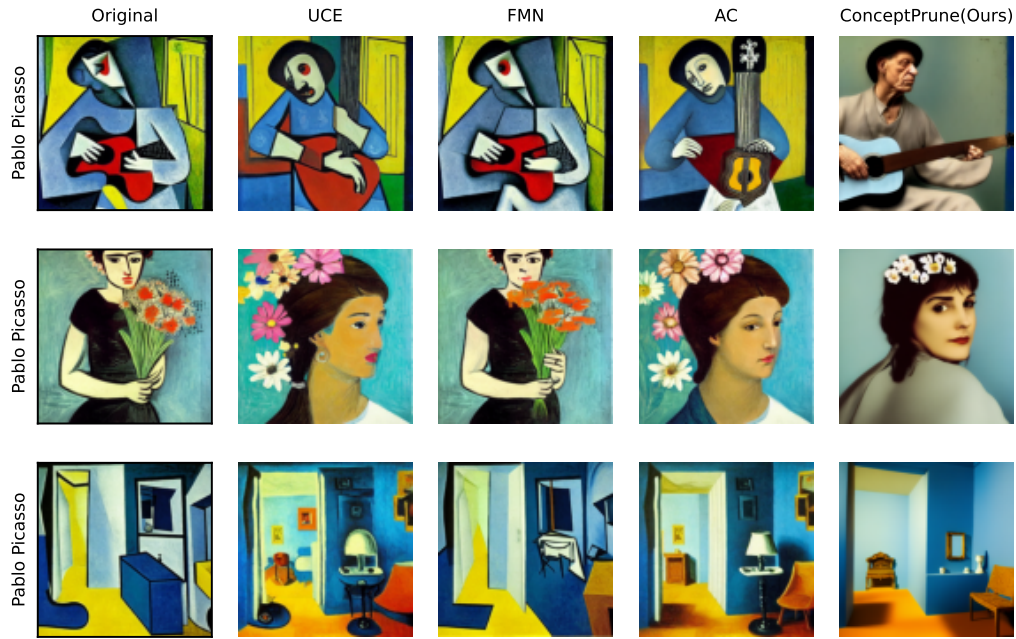


Figure 8: Qualitative results for erasing artist - *Pablo Picasso*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist’s style.

996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

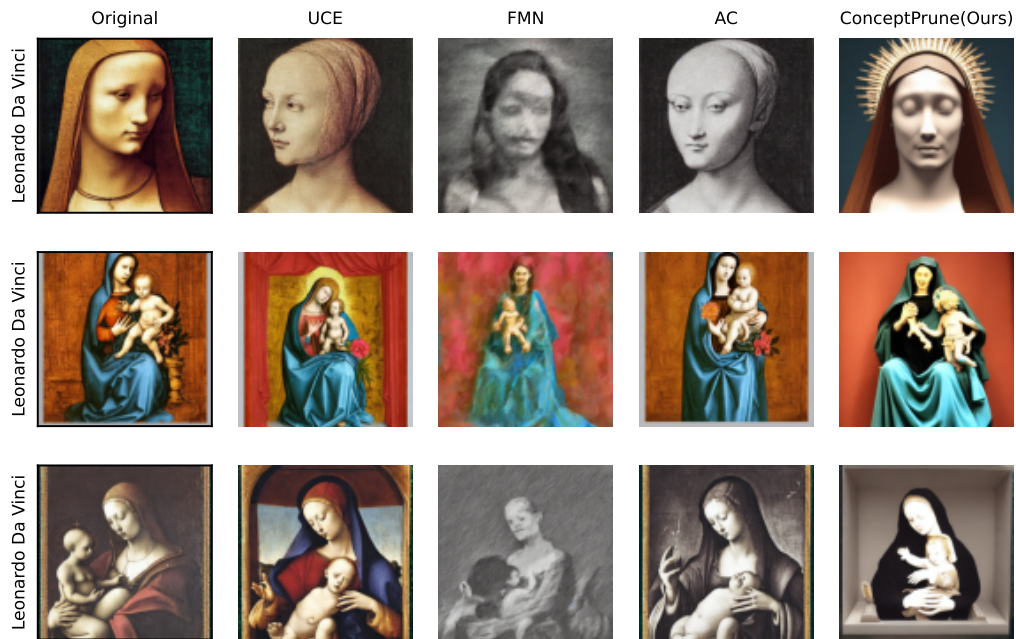


Figure 9: Qualitative results for erasing artist - *Leonardo da Vinci*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist’s style.

Table 10: Details on hyper-parameters, sparsity level and \hat{t} for concepts considered in our experiments.

Global Concept	Concept	Sparsity Level $k\%$	\hat{t}
Art Styles	Van Gogh	2.0	10
	Monet	2.0	10
	Leonardo Da Vinci	2.0	10
	Salvador Dali	2.0	10
	Pablo Picasso	2.0	10
Nudity	naked	1.0	9
Object Erasure	ImageNette classes	2.0	10
Gender change	Male to Female	5.0	20
	Female to Male	5.0	20

Table 11: Extension of Table 2 for Artist Style removal in the main paper. We report CLIP Similarity and CLIP Accuracy for 5 artists.

Artist	Metric	ESD	UCE	FMN	CA	SA	MACE	Receler	AdvUnlearn	ConceptPrune
Van Gogh	CLIP Similarity	33.1	34.3	26.6	32.9	24.5	27.8	30.1	28.5	29.2
	CLIP Accuracy (%)	39.0	36.0	96.0	58.0	96.0	82.5	79.8	69.0	84.0
Claude Monet	CLIP Similarity	32.9	33.6	23.2	33.1	25.6	24.5	23.9	25.0	23.6
	CLIP Accuracy (%)	57.0	56.0	98.0	68.0	94.9	95.7	98.2	97.6	100
Pablo Picasso	CLIP Similarity	33.5	32.9	33.0	31.3	30.9	28.9	29.3	26.1	25.3
	CLIP Accuracy (%)	58.0	56.0	58.0	78.0	72.0	75.6	78.4	82.7	100
Leonardo Da Vinci	CLIP Similarity	30.8	31.5	25.1	31.6	24.5	27.1	25.7	26.3	26.5
	CLIP Accuracy (%)	66.0	64.0	62.0	56.0	87.6	88.1	73.2	65.3	94.0
Salvador Dali	CLIP Similarity	39.9	31.6	33.6	32.8	30.1	32.7	33.1	29.9	29.8
	CLIP Accuracy (%)	26.0	8.0	98.0	66.0	83.9	85.2	80.3	95.2	92.0

A.6 CONCEPT INVERSION

We present the results of baselines considered in the paper in Table 13, which shows that ConceptPrune offers significantly greater adversarial robustness against CI.

A.7 ARE THERE SPECIFIC NEURONS RESPONSIBLE FOR GENERATING GENDER?

It is widely acknowledged that image-generation models harbor societal and gender biases [Luccioni et al., 2023]. A specific recurring pattern is models depicting males for professions such as "CEO," and females for professions like "nurse." Concept editing methods like UCE [Gandikota et al., 2023b] and MEMIT [Orgad et al., 2023] have addressed these issues by debiasing models to ensure an equal representation of males and females across all professions. However, Gemini [et al, 2024] recently faced criticism for controversies stemming from over-debiasing models, resulting in the generation of factually or historically incorrect information⁴. This occurs because while debiasing may show a range of people for some cases, it fails to appropriately handle cases where such variation is not applicable.

To address this, we believe that gender choice in diffusion models should be precisely controllable, e.g., under the guidance of expert ethics committees. To explore, this we illustrate controlled Gender

⁴Our intention is not to defame. We only use this incident to motivate controlled gender reversal.

Model	Monet	Salvador Dali	Pablo Picasso	Da Vinci	Average
UCE	32.4	30.3	28.8	29.8	30.3
AC	31.6	28.9	26.7	28.4	28.9
ConceptPrune (Ours)	30.9	31.2	29.9	28.7	30.2

Table 12: We erase 'Van Gogh' style from the model and report CLIP similarity (\uparrow) on surrounding artist styles. Higher CLIP similarities indicate better preservation of surrounding artist styles.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

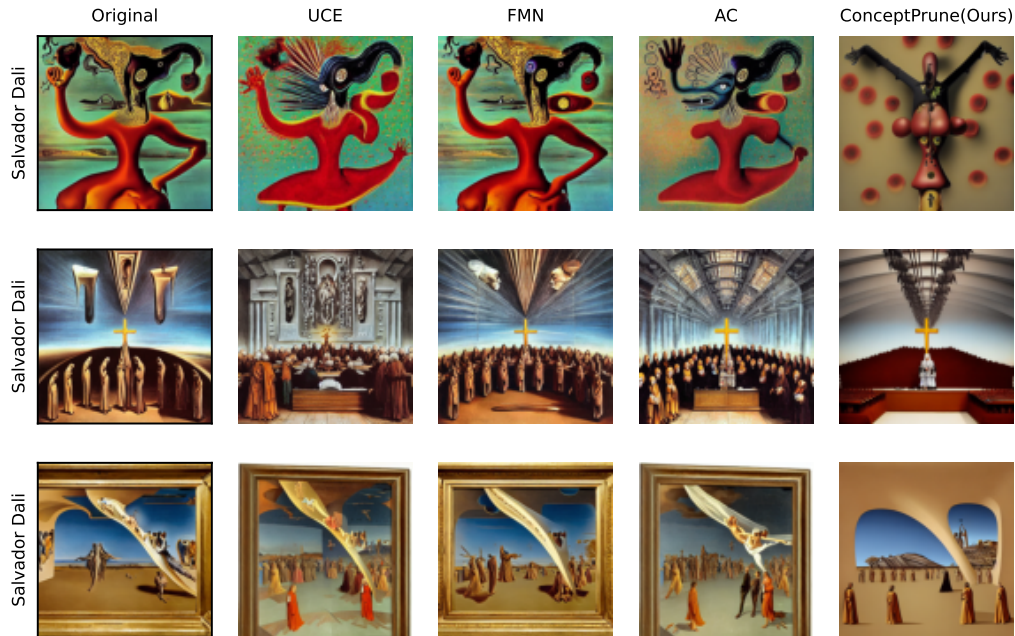


Figure 10: Qualitative results for erasing artist - *Salavdor Dali*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist’s style.

	ESD	FMN	UCE	CA	Neg-Prompt	SLD-Med	ConceptPrune (Ours)
Tench	59.7	60.6	20.6	29.4	72.6	75.4	0.0
Church	87.4	0.0	82.2	72.6	78.4	72.0	11.0
Parachute	94.2	93.4	94.2	92.4	77.2	95.8	0.0
Garbage Truck	57.0	69.6	89.6	79.4	84.6	94.8	6.8
Average	74.5	55.9	71.7	68.5	78.2	84.5	4.5

Table 13: Top-1 classification accuracy (\downarrow) under CI [Pham et al., 2024] for 4 Imagenette classes.

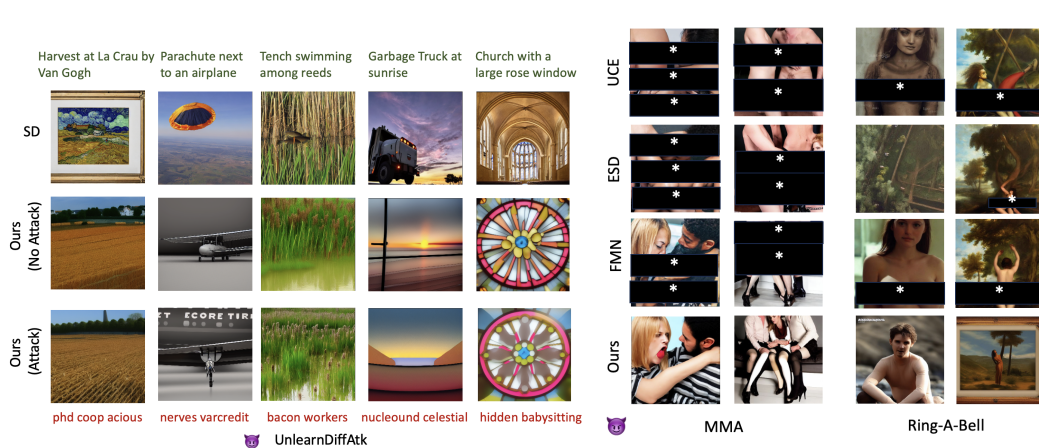


Figure 11: Qualitative results of the failure cases of adversarial attacks demonstrating the robustness of ConceptPrune to both white-box and black-box adversaries. *Left*: Top, middle, and bottom rows correspond to images generated by original SD, ConceptPrune without attack, and ConceptPrune under white-box UnlearnDiffAtk attack respectively. *Right*: Qualitative results of black-box attacks MMA [Yang et al., 2023] and Ring-A-Bell [Tsai et al., 2024] along with quantitative results in 3 show that ConceptPrune maintains its content moderation abilities even under attacks.

SD 2.0	Van Gogh	Monet	Salvador Dali	Pablo Picasso	Da Vinci	SD-XL	Van Gogh	Monet	Salvador Dali	Pablo Picasso	Da Vinci
UCE	32.5	25.6	31.8	25.8	26.9	UCE	31.4	29.3	28.4	27.8	27.6
ConceptPrune	30.2	23.7	28.8	24.1	25.7	ConceptPrune	29.4	27.8	29.0	24.7	26.7

Table 14: CLIP similarity (\downarrow) for artist erasure experiments with SD-v2.0 (left) and SD-XL (right). Our ConceptPrune can effectively erase artist styles.

Reversal⁵. We discover a set of “male” neurons via concept prompts \mathcal{P}^* like {a man, a boy}, vs reference prompts \mathcal{P} like {a woman, a girl} and vice-versa. Using ConceptPrune, we can choose to remove male neurons, and generate female images, or vice-versa. This allows direct control of gender for any future prompt, via simple choice of mask. We evaluate our model across 35 professions in the Winobias dataset [Zhao et al., 2018] and report the *success rate* at which the gender of the individual as classified by CLIP was reversed by ConceptPrune as compared to pre-trained SD. Qualitative results for controlled gender reversal are presented in Figure 3 (Left). We observed that our model has a *success rate* of $87 \pm 12\%$ with more failure cases like erasing the person from the image arising from highly male or female-biased professions like Carpenter, Secretary, etc. In this paper, we do not propose ConceptPrune as a practical solution for mitigating gender bias. Instead, our primary objective is to emphasize the compelling discovery of a distinct set of gender-specific neurons within the model.

⁵We exclude non-binary genders to ensure a clear evaluation of gender reversal success rates.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

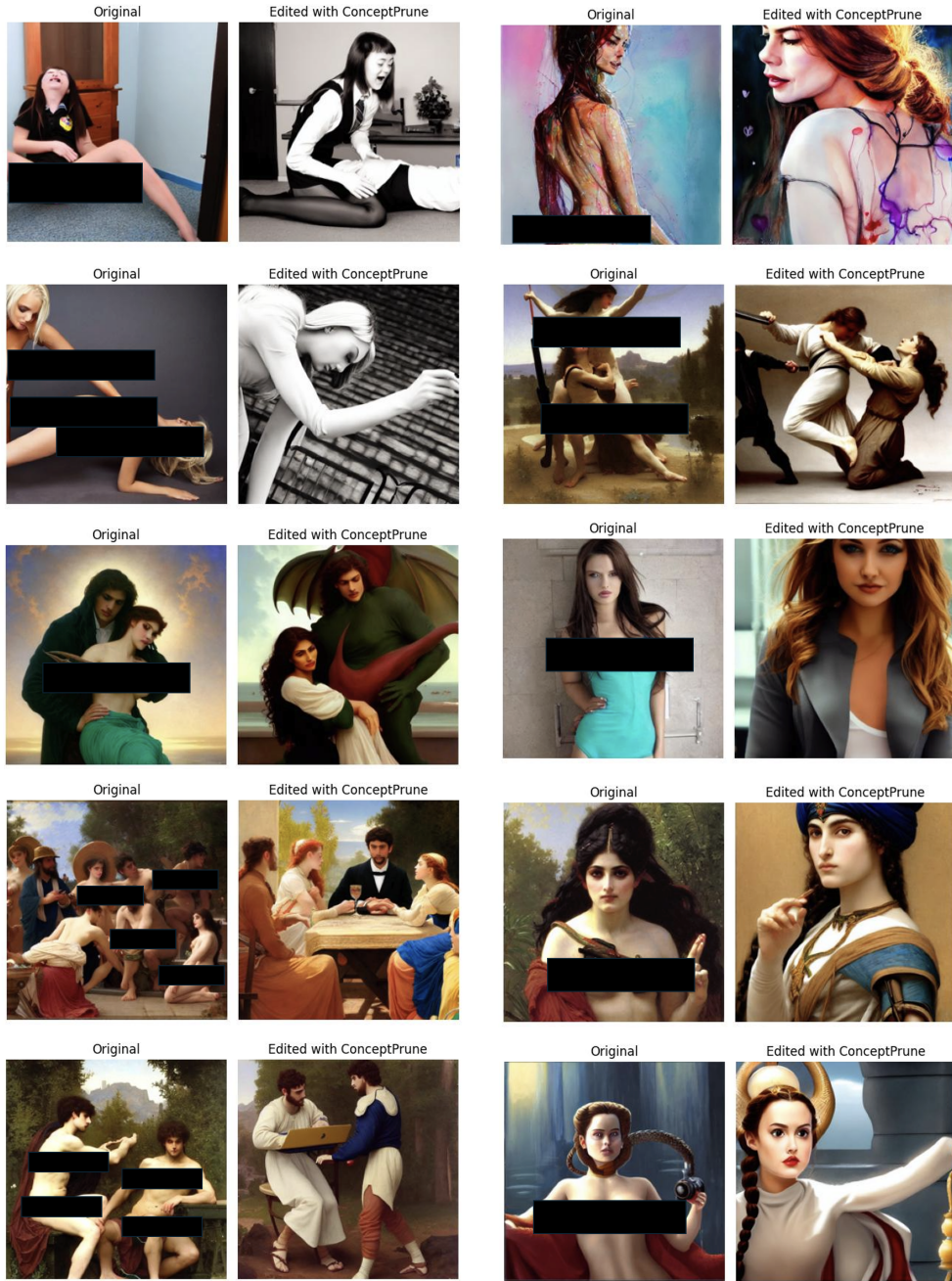


Figure 13: Qualitative results for Nudity Erasure. We omit the prompts for safety. Images marked as "Original" correspond to images generated by pre-trained Stable Diffusion. Sensitive parts have been blacked out by the authors for the purpose of publication. We observe that ConceptPrune erases nudity while preserving other details and quality of the image.

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295



Figure 12: Qualitative visualizations of controlled Gender Reversal using ConceptPrune. $M \rightarrow F$ and $F \leftarrow M$ indicate the removal of “male” generating and “female” generating neurons respectively. In most cases ConceptPrune succeeds in reversing the gender of the individual.

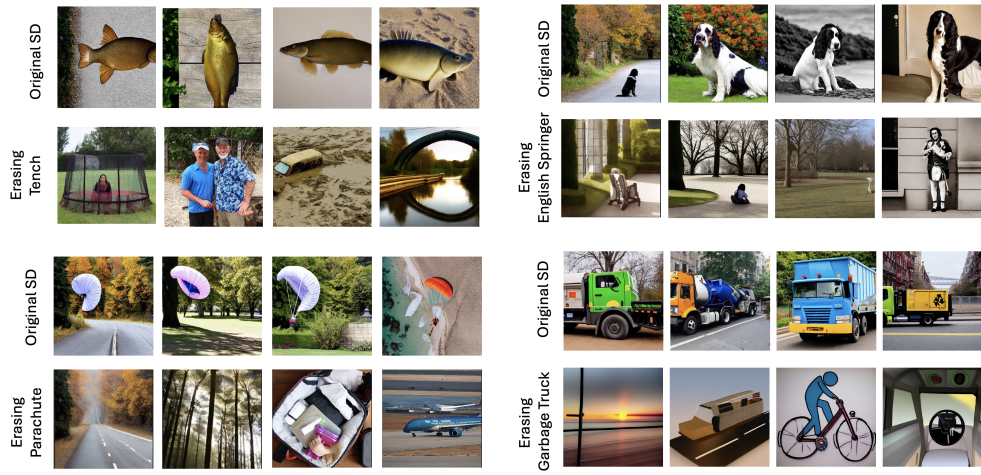


Figure 14: Qualitative results for Object Erasure