

---

# ConceptPrune: Concept Editing in Diffusion Models via Skilled Neuron Pruning

---

Ruchika Chavhan<sup>1</sup>, Da Li<sup>2</sup>, Timothy Hospedales<sup>1,2</sup>

<sup>1</sup>University of Edinburgh,

<sup>2</sup>Samsung AI Research Centre, Cambridge

## 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, object erasure, and gender debiasing 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. Our code is available at <https://github.com/ruchikachavhan/concept-prune.git>

## 1 Introduction

In recent years, text-to-image generation has witnessed significant advances driven by the development and adoption of diffusion models (DMs) [24, 43, 45, 46, 35, 60, 33, 39] across industries and real-world scenarios. However, this swift advancement presents a substantial risk. Diffusion models can threaten artists' livelihoods through style replication [11], generate convincing deepfakes and NSFW content [40, 14], and perpetuate societal biases [32]. 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 [48], COYO [4], and CC12M [5], 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 [26, 18, 19, 58, 36] and Model Unlearning (MU) [23, 65, 30, 56, 12], 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 [26, 18] and token remapping [58, 19]. 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 resteeering [58], adversarial attacks based on textual inversion [61, 38, 57, 53] have demonstrated the

---

Correspondence to: {ruchika.chavhan, t.hospedales}@ed.ac.uk

ability to circumvent concept erasure methods [18, 19, 58] that were previously believed to be robust with a near-perfect success rate.

In this paper, we introduce *ConceptPrune*, an entirely training-free method for concept editing that, for the first time, tackles knowledge editing in diffusion models through the lens of pruning. Leveraging recently introduced pruning heuristics [52], we identify regions or neurons in feed-forward layers of diffusion models that strongly activate in the presence of a concept, and denote them as *skilled neurons*. Subsequently, concept removal can be achieved by simply pruning or *zeroing* out these skilled regions. We demonstrate that *ConceptPrune* provides a rapid, efficient, and unified solution for erasing undesired concepts, including various artist styles, nudity, undesired objects, and gender biases. Notably, it maintains the outstanding image-generation prowess of pre-trained models while remaining resilient to adversarial attacks.

## 2 Related Work

**Diffusion model concept editing:** Most concept Editing works [26, 18, 19, 36, 58] within diffusion models aim to eliminate target concepts by aligning the model’s output with that of a reference prompt, whose concept we wish to retain. For example, to remove the concept ‘nudity’, the target prompt can be formulated as “*a photo of a naked person*” while the reference prompt can be “*a photo of a person*”. Then the target concept “nudity” can be removed by minimizing certain metrics between denoised predictions of target and reference prompts [26], utilizing score-based composition as unsupervised training data [18], or employing attention re-steering to reduce cross-attention weights for the target prompts [58]. Unlike approaches relying on latent representations, UCE [19] and MEMIT [36] operate on token rewriting, adjusting attention module parameters in the UNet to align token embeddings corresponding to the target prompt with the reference prompt using a closed-form solution.

**Diffusion model unlearning (MU):** MU [65, 23, 12, 56, 30] operates with two separate datasets: the forgetting dataset and the retention dataset. The model is fine-tuned such that information from the forgetting dataset is erased while knowledge corresponding to the remaining data remains intact. There are different ways to achieve this dual objective optimization: a first-order dual problem formulation [56], generative replay on the retention dataset to ensure consistent retention of the dataset [23], and fine-tuning via saliency masks that retain the reference concept while disregarding the target concept [12]. While these methods have shown remarkable efficacy in unlearning multiple concepts, they are usually computationally expensive, especially for large-scale models.

In contrast to the existing Concept Editing and MU methods, our method operates on training-free neuron identification and pruning of critical regions that are responsible for generating undesired behaviors. While our method does not necessitate any compute-intensive fine-tuning of parameters, it is directly aligned with the Concept Editing line of work as we demonstrate that denoised prediction matching is possible through the selective removal of neurons in the weight space.

**Language model skilled neuron identification:** Previous works [54, 51, 8, 6, 9, 1] present strong evidence that activation of specific neurons in feed-forward networks in transformers show high correlation with task labels, with perturbations to these neurons impacting task performance. Modular components within pre-trained transformers were identified by leveraging the inherent sparsity in neurons, as shown in [62]. Further, [63] demonstrates that these modules are specialized in distinct functions. In this work, we aim to identify neurons 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 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 [27, 31, 20, 15, 3] aims to reduce model size either by eliminating parameters and substructures from networks [29, 16] or by masking parameters guided by a score function [16, 17, 52, 28]. This study primarily focuses on the latter approach. Exploration of diffusion model pruning is limited, although one study [13] introduces structural pruning by accumulating gradient-based importance scores across a chosen subset of denoising time steps.

One exploratory study [55] delves into safety-aligned large language models (LLMs) that [37] possess the ability to inhibit responses to harmful prompts. They leverage heuristics from diverse pruning methods [52, 28] to pinpoint the regions that deny harmful responses to triggering prompts. Further,

they illustrate that these regions lie within a compact zone in the weight space and their removal poses a huge risk to safety alignment in LLMs. In contrast, we derive insights from model pruning heuristics [52] to pinpoint critical regions in the weight space accountable for unsafe behaviors already learned by the pre-trained model and subsequently unlearn them permanently through model pruning.

### 3 Preliminaries

**(Latent) diffusion models:** Diffusion models (DMs) [24, 50] 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) [41, 59] are increasingly used as the first choice of DMs as they accelerate the above process by operating in a latent space, denoted as  $\mathbf{z}$ , 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 certain concepts activate specific neurons in a neural network [34, 54], 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 [43], characterized by a UNet [44] denoted as  $f_\theta$ . 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 [49]. Prior research in concept editing, such as [18] and [58], 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 [51, 54, 63, 9, 6], 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 [49] 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 [22]. 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 [52] 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  $\ell_2$  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  $\ell_2$  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 [52]. 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* [52]. 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 [52], 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 [13] 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  $\vee$  and  $\neg$  denote the logical OR and NOT operators. All the weights of the pre-trained diffusion model  $f_\theta$  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.

## 5 Experiments

**Experimental details:** We work with Stable Diffusion-v1.5 (SD), which includes 16 FFN layers that serve as candidates for skilled neuron discovery and pruning. We begin by formulating the calibration sets  $\mathcal{P}^*$  and  $\mathcal{P}$  that are 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 prompts and the exact structure of the sentences for different concepts is provided in Table 6 in the Appendix. To calculate neuron activations, we run the model for 50 denoising iterations and fix the seed before every forward pass to ensure the same initializations for both reference and target concept prompts. As discussed in Section 4.1, we select two hyperparameters sparsity level  $k\%$  and  $\hat{t}$  for aggregating skilled neurons over time steps. The values of sparsity levels  $k\%$  and hyperparameter  $\hat{t}$  chosen for each concept are detailed in Table 7 in the Appendix. Interestingly, our experiments show that  $\hat{t} = 10$  is sufficient for removing concepts while retaining objects for most cases, suggesting that low-level features like style and objects are generated early in the denoising process, followed by the addition of fine-grained details.

**Baselines:** We consider the following concept editing works as our closest competitors: UCE [19], ESD [18], Forget-Me-Not (FMN) [58], and Concept Ablation (CA) [26]. While ESD, UCE, and FMN experiment with erasing artist styles, objects, and nudity, CA does not evaluate their method on nudity and the same objects. Therefore, we include a baseline only if their method has been evaluated for that concept and is reproducible from their source code <sup>1</sup>.

### 5.1 Erasing Artistic Styles

We consider five artists — *Van Gogh*, *Claude Monet*, *Pablo Picasso*, *Leonardo Da Vinci*, and *Salvador Dali*. To measure the efficacy of concept removal, we created a dataset of 50 prompts for each artist using ChatGPT, consisting of the names of their paintings followed by the name of the artist. To measure the efficacy of concept removal, we report two metrics: the *CLIP Similarity*, which measures the similarity between the generated image and the prompt, and a stricter *CLIP score* that penalizes a model when the similarity between the image generated by the concept-editing and the prompt is greater than the similarity between the image generated by the pre-trained SD and prompt. Lower

<sup>1</sup>We reproduced CA to remove nudity and object classes from ImageNette but performance was very poor.

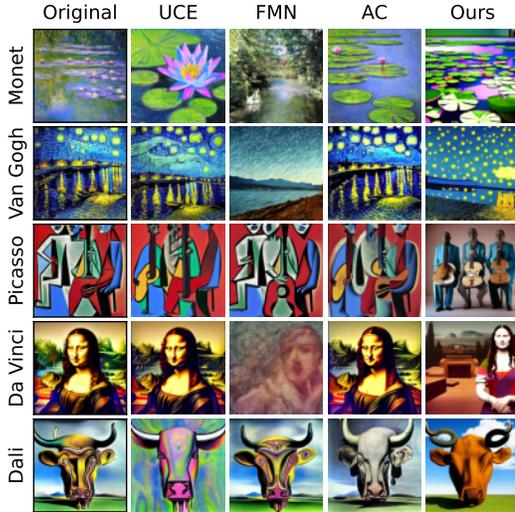


Figure 1: Qualitative results of artist erasure. ConceptPrune demonstrates stronger erasing while generating high-quality, realistic-looking images.

Table 1: 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 8. Our ConceptPrune can effectively erase artist styles without compromising the model’s performance on unrelated concepts.

	Artist erasure		COCO	
	Similarity ↓	Score ↑	FID ↓	Similarity ↑
Original SD	42.1	23.0	<b>14.5</b>	<b>31.3</b>
ESD [18]	34.1	49.2	15.9	30.7
UCE [19]	32.8	44.0	15.7	<b>31.3</b>
FMN [58]	28.4	82.4	20.9	29.8
CA [26]	32.4	65.2	17.5	31.3
ConceptPrune	<b>26.9</b>	<b>94.0</b>	16.9	29.9

values of *CLIP Similarity* and higher values of *CLIP Score* indicate better concept removal. We also evaluate the fidelity of general purpose image generation by measuring FID and *CLIP Similarity* on the COCO30k dataset. From the quantitative results presented in Table 1, we demonstrate that our method outperforms other baselines in artist style removal while effectively retaining unrelated concepts, as indicated by the low FID score. In Figure 1, we present some qualitative examples that demonstrate the strong erasing capabilities of ConceptPrune with high-quality realistic output images. More qualitative results are presented in Section A.2 in the appendix.

## 5.2 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) [47] and Stable Diffusion 2.0 [42], which have been fine-tuned on a filtered subset of LAION without explicit images. We use the Inappropriate Prompts Dataset (I2P) [47], which consists of 4703 prompts featuring various inappropriate concepts. Nudity detectors [2] 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 11 in the appendix.

## 5.3 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 ImageNet classes [25], a subset of ImageNet [7] 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 [21]. Table 3 shows that ConceptPrune has superior erasure performance on average while effectively minimizing interference on non-targeted classes. More results of object erasure are provided in Figure 12 in the appendix.

**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

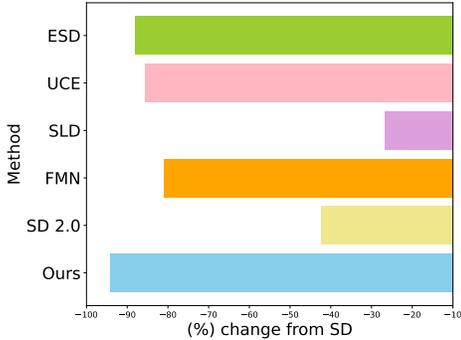


Figure 2: Explicit Content Erasure. The percentage reduction in nudity content from I2P prompts, compared to the original SD model ConceptPrune (SD1.5) decreases the number of explicit images by 94.1%, outperforming competitors as well as SD2.0.

Table 2: ConceptPrune demonstrates robustness to adversarial attacks. Unlearning methods evaluated against three adversarial attacks. Black-box (Ring-A-Bell[53], and MMA[57]) performance is quantified by percentage reduction in nude samples compared to SD. White-box UnlearnDiffAtk [61] performance measures the attack success rate (ASR).

	Ring-A-Bell $\uparrow$	MMA $\uparrow$	UnlearnDiffAtk $\downarrow$
ESD [18]	52.8	87.3	76.1
UCE [19]	67.6	63.3	93.2
SLD [47]	2.80	25.5	82.4
FMN [58]	5.60	53.6	97.9
SDv2 [43]	1.80	26.8	73.8
Ours	<b>85.2</b>	<b>95.6</b>	<b>64.8</b>

Table 3: 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 [18]	UCE [19]	FMN[58]	ConceptPrune	ESD [18]	UCE [19]	FMN [58]	ConceptPrune
Church	54.2	8.4	2.0	<b>1.9</b>	<b>80.2</b>	77.5	57.8	74.5
English Springer	6.2	0.2	1.9	<b>0.0</b>	62.6	78.9	73.5	<b>93.7</b>
Golf ball	5.8	<b>0.8</b>	13.7	6.9	65.6	79.0	82.8	<b>98.6</b>
Gas Pump	8.6	<b>0.0</b>	7.9	<b>0.0</b>	66.5	<b>80.7</b>	79.0	79.1
Tench	9.6	<b>0.0</b>	5.7	<b>0.0</b>	66.6	79.3	78.4	<b>90.1</b>
Parachute	23.8	<b>1.4</b>	8.3	6.9	65.4	77.4	<b>98.2</b>	72.8
Cassette Player	0.6	<b>0.0</b>	1.0	1.9	64.5	<b>90.3</b>	68.7	82.5
Chain Saw	6.0	<b>0.0</b>	0.1	<b>0.0</b>	71.6	<b>80.2</b>	78.4	77.8
French Horn	0.4	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	77.0	80.1	78.3	<b>81.1</b>
Garbage Truck	10.4	14.8	0.1	<b>0.0</b>	51.5	<b>78.7</b>	74.9	75.2
Average	12.5	2.7	4.1	<b>1.8</b>	66.9	<b>80.2</b>	77.5	<b>80.4</b>

erasing by taking the union of skilled neurons across different concepts. We direct the reader to Appendix A.3 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 5 shows that ConceptPrune demonstrates comparable erasing performance while excelling at retaining unrelated concepts.

#### 5.4 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 introduced in [61] 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. We evaluate ConceptPrune under UnlearnDiffAtk 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 [61], 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 [2]. Table 4 illustrates that for artist style and object erasure, ConceptPrune renders the attack unsuccessful, achieving a 0% ASR in two instances, in contrast to the perfect success rates seen for baselines like UCE and FMN. Table 2 shows that UCE, ESD,

Table 4: ConceptPrune is substantially more robust to adversarial attacks aimed at eliciting erased concepts. Attack Success Ratio (ASR %, ↓) of UnlearnDiffAtk [61] adversarial prompts for Van Gogh’s painting style and 4 classes of the Imagenette dataset.

	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 [18]	32.0	76.0	54.0	36.0	24.0	60.0
UCE [19]	94.0	100.0	43.0	22.0	38.0	68.0
FMN [58]	56.0	90.0	100.0	100.0	98.0	96.0
CA [26]	77.0	92.0	–	–	–	–
ConceptPrune (Ours)	<b>2.0</b>	<b>24.0</b>	<b>34.0</b>	<b>10.4</b>	<b>0.0</b>	<b>22.2</b>

Table 5: Quantitative results for multi-object erasure. We report Accuracy on erased classes and FID on COCO30k and CLIP similarity on COCO30k. ConceptPrune is comparable to UCE at erasing multiple objects and outperforms UCE in retaining image generation capabilities.

	COCO FID	CLIP score	Accuracy on erased classes
UCE [19]	17.7	31.0	4%
ConceptPrune	17.5	29.9	7%

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). We present more qualitative analysis in Figure 9 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 [57] 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 [53] 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 [57, 53] and report the percentage reduction in number of images for which nudity is generated as compared to pre-trained SD. Results in Table 2 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 9 in the appendix.

## 5.5 Gender Reversal

It is widely acknowledged that image-generation models harbor societal and gender biases [32]. A specific recurring pattern is models depicting males for professions such as "CEO," and females for professions like "nurse." Concept editing methods like UCE [19] and MEMIT [36] have addressed these issues by debiasing models to ensure an equal representation of males and females across all professions. However, Gemini [10] recently faced criticism for controversies stemming from over-debiasing models, resulting in the generation of factually or historically incorrect information<sup>2</sup>. 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 Reversal<sup>3</sup>. 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

<sup>2</sup>Our intention is not to defame. We only use this incident to motivate controlled gender reversal.

<sup>3</sup>We exclude non-binary genders to ensure a clear evaluation of gender reversal success rates.

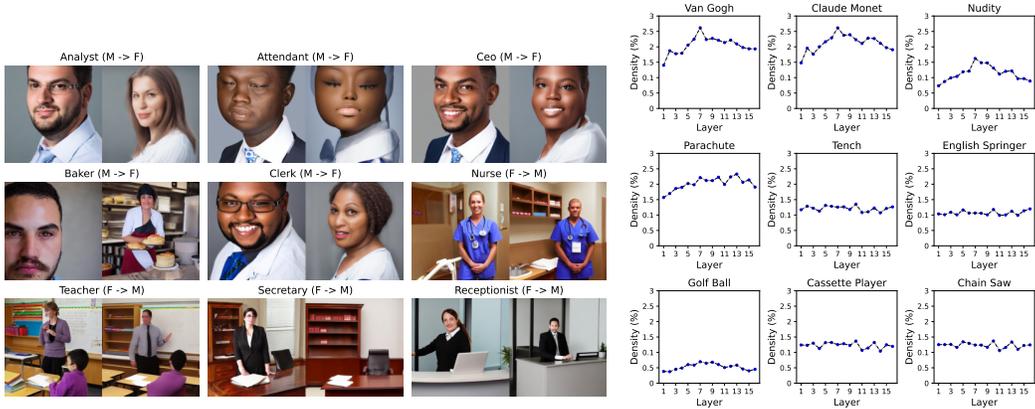


Figure 3: *Left*: Qualitative visualizations of controlled Gender Reversal using ConceptPrune.  $M \rightarrow F$  and  $F \rightarrow M$  indicate the removal of “male” generating and “female” generating neurons respectively. In most cases ConceptPrune succeeds in reversing the gender of the individual. *Right*: Skilled neurons are localized to a very compact subspace, between 1% to 3% of FFN parameters.

in the Winobias dataset [64] 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.

## 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 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. We present more interesting analysis on the disentangled nature of skilled neurons in Sec A.4 in the appendix.

## 6 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. In our experiments with controlled gender reversal, we observed that while ConceptPrune successfully reverses gender in majority of instances, it sometimes also removes the person from the image. 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.

## 7 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

- [1] Omer Antverg and Yonatan Belinkov. On the pitfalls of analyzing individual neurons in language models. *ICLR*, 2022.
- [2] Praneeth Bedapudi. Nudenet: Neural nets for nudity detection and censoring. 2022.
- [3] Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Guttery. What is the state of neural network pruning? *MLSys*, 2020.
- [4] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. 2022.
- [5] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021.
- [6] Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James Glass. What is one grain of sand in the desert? analyzing individual neurons in deep nlp models. *AAAI*, 2018.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: a large-scale hierarchical image database. 2009.
- [8] Nadir Durrani, Fahim Dalvi, and Hassan Sajjad. Discovering salient neurons in deep nlp models. *JMLR*, 2023.
- [9] Nadir Durrani, Hassan Sajjad, Fahim Dalvi, and Yonatan Belinkov. Analyzing individual neurons in pre-trained language models. *EMNLP*, 2020.
- [10] Gemini Team et al. Gemini: A family of highly capable multimodal models. *arXiv*, 2024.
- [11] Sarah Andersen. et al v. Stability AI Ltd. et al. Case no.3:2023cv00201. us district court for the northern district of california., 2023.
- [12] Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. *ICLR*, 2024.
- [13] Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. *NeurIPS*, 2023.
- [14] Camera Forensics. The dark reality of stable diffusion. 2024.
- [15] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *ICLR*, 2019.
- [16] Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. *ICML*, 2023.
- [17] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *ICLR*, 2023.
- [18] Rohit Gandikota, Joanna Materzyńska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *ICCV*, 2023.
- [19] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified concept editing in diffusion models. *WACV*, 2023.
- [20] Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks. *NeurIPS*, 2015.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CVPR*, 2015.
- [22] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv*, 2023.

- [23] Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep generative models. In *NeurIPS*, 2023.
- [24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020.
- [25] Jeremy Howard and Sylvain Gugger. Fastai: A layered api for deep learning. *Information*, 11, 2020.
- [26] Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu. Ablating concepts in text-to-image diffusion models. In *ICCV*, 2023.
- [27] Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In *Advances in Neural Information Processing Systems*, 1989.
- [28] Namhoon Lee, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Snip: Single-shot network pruning based on connection sensitivity. *ICLR*, 2019.
- [29] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *ICLR*, 2017.
- [30] Zhili Liu, Kai Chen, Yifan Zhang, Jianhua Han, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan Yeung, and James Kwok. Implicit concept removal of diffusion models. *arXiv*, 2024.
- [31] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network pruning. In *ICLR*, 2019.
- [32] Sasha Luccioni, Christopher Akiki, Margaret Mitchell, and Yacine Jernite. Stable bias: Evaluating societal representations in diffusion models. *NeurIPS*, 2023.
- [33] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv*, 2023.
- [34] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In *CVPR*, 2015.
- [35] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *ICML*, 2022.
- [36] Hadas Orgad, Bahjat Kawar, and Yonatan Belinkov. Editing implicit assumptions in text-to-image diffusion models. *ICCV*, 2023.
- [37] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. *NeurIPS*, 2022.
- [38] Minh Pham, Kelly O. Marshall, and Chinmay Hegde. Circumventing concept erasure methods for text-to-image generative models. *ICLR*, 2024.
- [39] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *ICLR*, 2024.
- [40] MIT Technology Review. Text-to-image ai models can be tricked into generating disturbing images. 2023.
- [41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.

- [43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *CVPR*, 2022.
- [44] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. *MICCAI*, 2015.
- [45] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. 2023.
- [46] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 2022.
- [47] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *CVPR*, 2023.
- [48] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models. In *NeurIPS*, 2022.
- [49] Noam Shazeer. Glu variants improve transformer. *arXiv*, 2020.
- [50] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021.
- [51] Xavier Suau, Luca Zappella, and Nicholas Apostoloff. Finding experts in transformer models. *arXiv preprint*, 2020.
- [52] Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. A simple and effective pruning approach for large language models. *ICLR*, 2024.
- [53] Yu-Lin Tsai, Chia-Yi Hsu, Chulin Xie, Chih-Hsun Lin, Jia-You Chen, Bo Li, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion models? *ICLR*, 2024.
- [54] Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. Finding skill neurons in pre-trained transformer-based language models. *EMNLP*, 2022.
- [55] Boyi Wei, Kaixuan Huang, Yangsibo Huang, Tinghao Xie, Xiangyu Qi, Mengzhou Xia, Prateek Mittal, Mengdi Wang, and Peter Henderson. Assessing the brittleness of safety alignment via pruning and low-rank modifications. *arXiv preprint arXiv:2402.05162*, 2024.
- [56] Jing Wu, Trung Le, Munawar Hayat, and Mehrtash Harandi. Erasediff: Erasing data influence in diffusion models. *arXiv preprint*, 2024.
- [57] Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Nan Xu, and Qiang Xu. Mma-diffusion: Multimodal attack on diffusion models. *CVPR*, 2024.
- [58] Eric Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not: Learning to forget in text-to-image diffusion models. *arXiv preprint arXiv:2211.08332*, 2023.
- [59] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, October 2023.
- [60] Tianyi Zhang, Zheng Wang, Jing Huang, Mohiuddin Muhammad Tasnim, and Wei Shi. A survey of diffusion based image generation models: Issues and their solutions. *IEEE PAMI*, 2023.
- [61] Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. *arXiv preprint*, 2023.

- [62] Zhengyan Zhang, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. Moefication: Transformer feed-forward layers are mixtures of experts. *ACL*, 2022.
- [63] Zhengyan Zhang, Zhiyuan Zeng, Yankai Lin, Chaojun Xiao, Xiaozhi Wang, Xu Han, Zhiyuan Liu, Ruobing Xie, Maosong Sun, and Jie Zhou. Emergent modularity in pre-trained transformers. *ACL*, 2023.
- [64] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in coreference resolution: Evaluation and debiasing methods. *NAACL*, 2018.
- [65] Mengnan Zhao, Lihe Zhang, Tianhang Zheng, Yuqiu Kong, and Baocai Yin. Separable multi-concept erasure from diffusion models. *arXiv*, 2024.

## A Appendix

### A.1 Details on prompts and hyper-parameters

Table 6: 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	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]
	Monet			
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 <sup>†</sup>
Object Erasure	parachute, gas pump	a <concept> in a <scene>	a <scene>	a <scene> ∈ [road, garden, beach, room, park, table, bag, tree, forest, street, shelter, chair]
	golf ball, cassette player			
english springer, tench				
chain saw, french horn				
church, garbage truck				
Object Erasure	church, garbage truck	a <concept> near a <place>	a <place>	<place> ∈ [road, park, beach, street, house, tree, forest, statue, car]
	Male to Female	a photo of a <male>	a photo of a <female>	<male> ∈ [man, boy, person, guy, father, son, husband, uncle]
Gender change	Female to Male	a photo of a <female>	a photo of a <male>	<female> ∈ [woman, girl, female, lady, mother, daughter, wife, aunt]

Table 7: 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

### A.2 Artist Style Erasure

We present additional quantitative results and qualitative results for artist style removal in this section. Please see Figure 4, 5, 6, 7, and 8 and Table 8.

### A.3 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  $\tilde{\mathbf{W}}_l^2$  as  $\mathbf{W}_l^2 \odot \left( \neg \left( \bigvee_{t=T, T-1, \dots, T-\hat{t}} \bigvee_{c \in \mathbb{C}} \mathbf{M}_c^{t,l} \right) \right)$ , where  $\bigvee$  and  $\neg$  denote the logical OR and NOT operators.

### A.4 Further Analysis

**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  $\mathbf{S}_i^l(\mathcal{P}^*)[i, j] < \mathbf{S}_i^l(\mathcal{P})[i, j]$  are hypothesised to distort the reference concept while retaining the target concept. Figure 10 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.

Table 8: Extension of Table 1 for Artist Style removal in the main paper. We report CLIP Similarity and CLIP Accuracy for 5 artists.

Artist	Metric	ESD	UCE	FMN	CA	ConceptPrune
Van Gogh	CLIP Similarity	33.1	34.3	<b>26.6</b>	32.9	29.2
	CLIP Accuracy (%)	39.0	36.0	<b>96.0</b>	58.0	84.0
Claude Monet	CLIP Similarity	32.9	33.6	<b>23.2</b>	33.1	23.6
	CLIP Accuracy (%)	57.0	56.0	98.0	68.0	<b>100</b>
Pablo Picasso	CLIP Similarity	33.5	32.9	33.0	31.3	<b>25.3</b>
	CLIP Accuracy (%)	58.0	56.0	58.0	78.0	<b>100</b>
Leonardo Da Vinci	CLIP Similarity	30.8	31.5	<b>25.1</b>	31.6	26.5
	CLIP Accuracy (%)	66.0	64.0	62.0	56.0	<b>94.0</b>
Salvador Dali	CLIP Similarity	39.9	31.6	33.6	32.8	<b>29.8</b>
	CLIP Accuracy (%)	26.0	8.0	<b>98.0</b>	66.0	92.0



Figure 4: Qualitative results for erasing artist - *Van Gogh*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist's style.

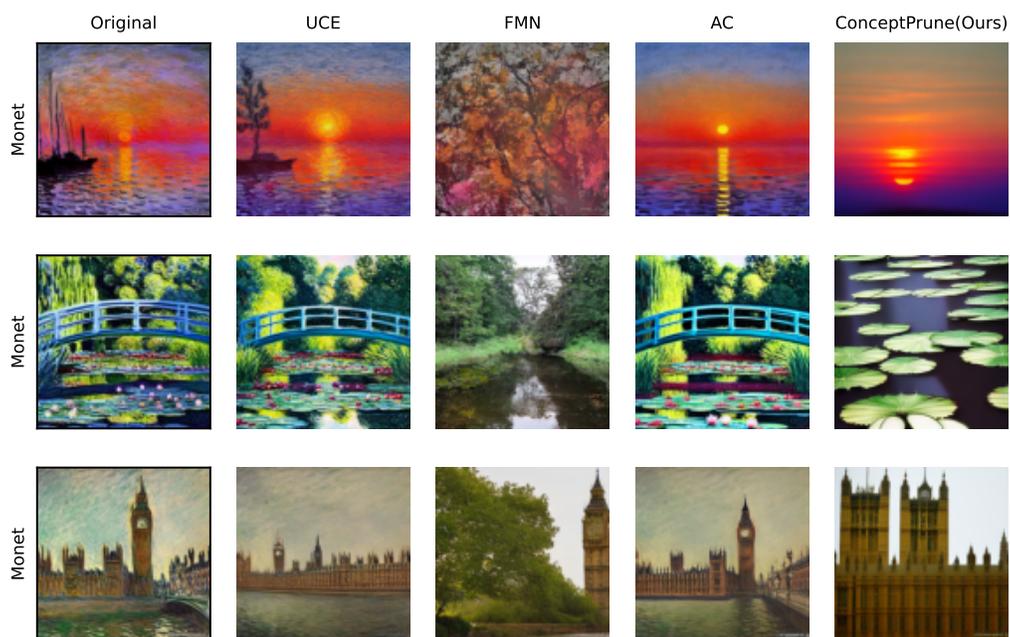


Figure 5: Qualitative results for erasing artist - *Monet*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist's style.

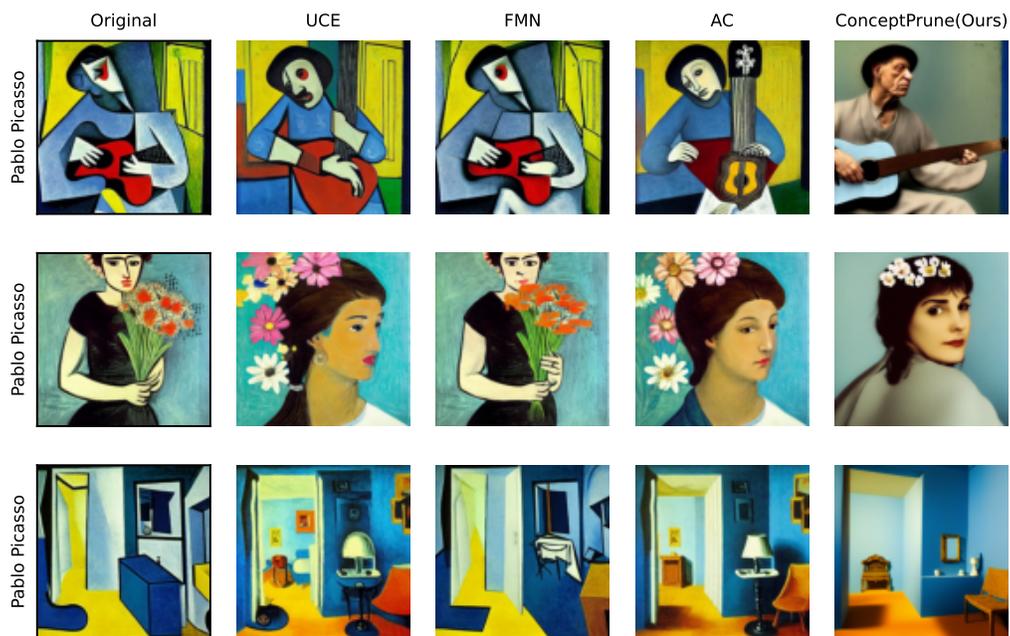


Figure 6: Qualitative results for erasing artist - *Pablo Picasso*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist's style.

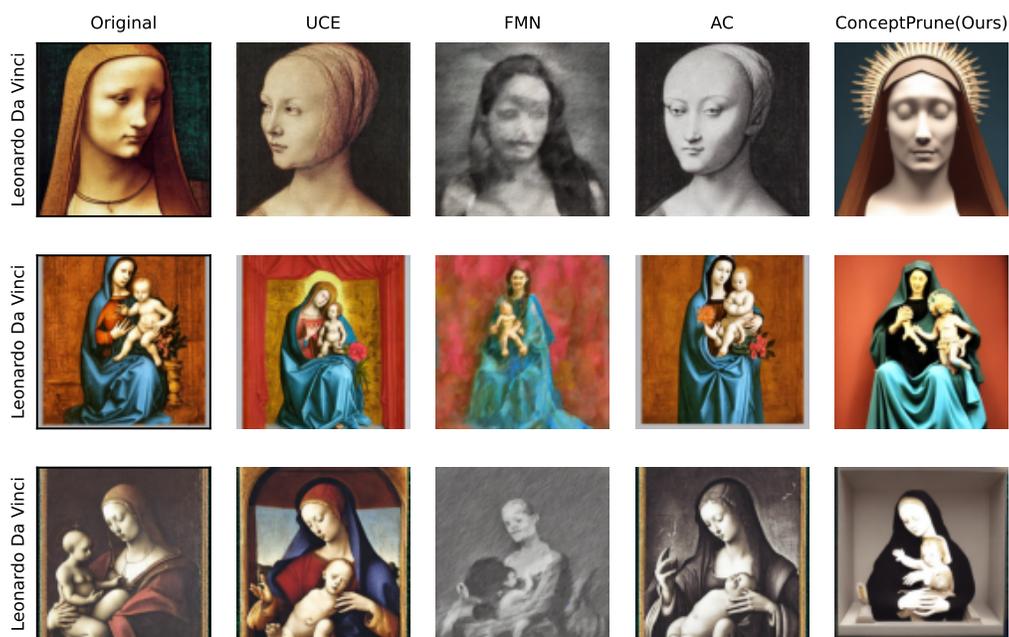


Figure 7: Qualitative results for erasing artist - *Leonardo da Vinci*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist's style.

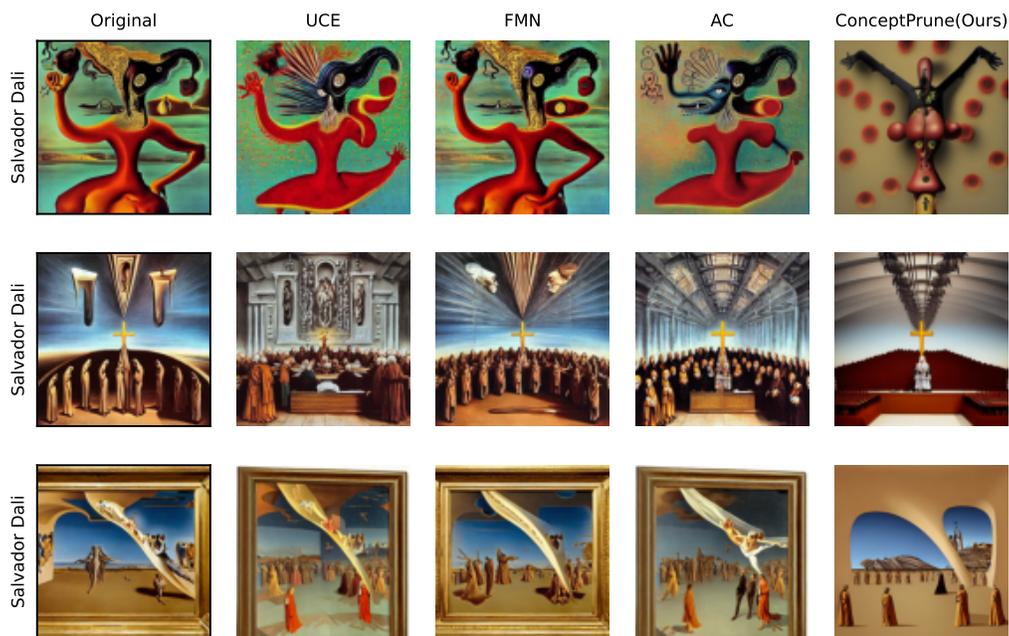


Figure 8: Qualitative results for erasing artist - *Salvador Dali*. ConceptPrune(Ours) generates high-quality realistic-looking images without the artist's style.



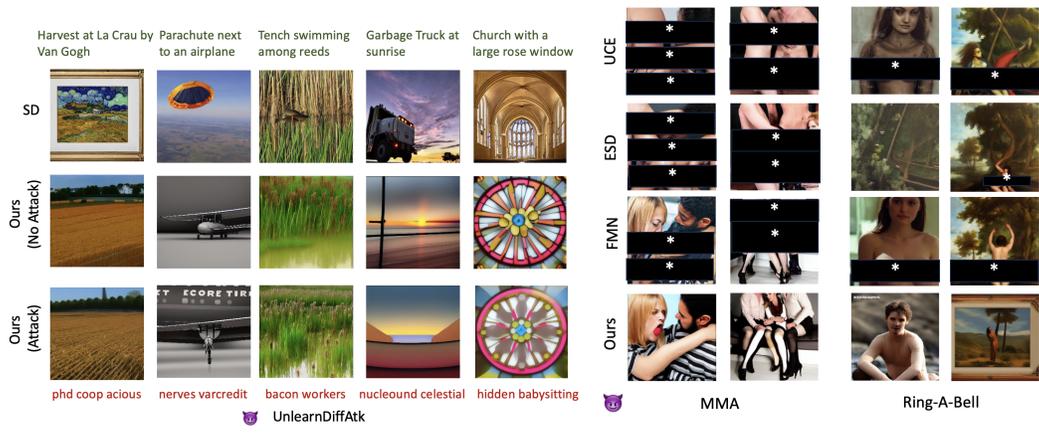


Figure 9: Qualitative results of the failure cases of adversarial attacks demonstrating the robustness of ConceptPrune to both white-box and black-box adversaries. *Left*: Top, middle, and bottom rows correspond to images generated by original SD, ConceptPrune without attack, and ConceptPrune under white-box UnlearnDiffAtk attack respectively. *Right*: Qualitative results of black-box attacks MMA[57] and Ring-A-Bell [53] along with quantitative results in 2 show that ConceptPrune maintains its content moderation abilities even under attacks.



Figure 12: Qualitative results for Object Erasure



Figure 10: 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 while preserving the "Van Gogh" style.