

TOWARDS ROBUST CONCEPT ERASURE IN DIFFUSION MODELS: UNLEARNING IDENTITY, NUDITY AND ARTISTIC STYLES

Anonymous authors

Paper under double-blind review

ABSTRACT

Diffusion models have achieved remarkable success in generative tasks across various domains. However, the increasing demand for content moderation and the removal of specific concepts from these models has introduced the challenge of *unlearning*. In this work, we present a suite of robust methodologies that significantly enhance the unlearning process by employing advanced loss functions within knowledge distillation frameworks. Specifically, we utilize the Cramer-Wold distance and Jensen-Shannon (JS) divergence to facilitate more efficient and versatile concept removal. Although current non-learning techniques are effective in certain scenarios, they are typically limited to specific categories such as identity, nudity, or artistic style. In contrast, our proposed methods demonstrate robust versatility, seamlessly adapting to and performing effectively across a wide range of concept erasure categories. Our approach outperforms existing techniques, achieving consistent results across different unlearning categories and showcasing its broad applicability. Through extensive experiments, we show that our method not only surpasses previous benchmarks but also addresses key limitations of current unlearning techniques, paving the way for more responsible use of text-to-image diffusion models.

1 INTRODUCTION

Diffusion models Ho et al. (2020); Dhariwal & Nichol (2021); Kawar et al. (2022); Gu et al. (2023) have advanced text-to-image generation, enabling the creation of high-quality visuals from diverse prompts. However, the extensive datasets used to train these models, often sourced indiscriminately from the Internet, pose significant ethical and safety challenges. It also presents a serious risk of misuse. Their ability to generate realistic images can be exploited to produce misleading or harmful content, such as deepfakes, disinformation, or unauthorized reproductions of copyrighted material. These concerns underscore the need for effective safeguards, including methods that can selectively *unlearn* or restrict certain concepts to prevent the misuse of these models.

Recently, there has been increasing interest in developing techniques to *unlearn* or *erase* specific concepts from diffusion models. In this direction, progress has been made in unlearning nudity, artistic, insignia, copyrighted images, and identity styles Kumari et al. (2023); Zhang et al. (2024a); Heng & Soh (2024); Gandikota et al. (2023); Kim et al. (2023); Golatkar et al. (2024); Lu et al. (2024). Although *concept erasure* is the primary objective of these methods, achieving consistent effectiveness in various unlearning scenarios remains a significant challenge. Identities of well-known figures, such as politicians or celebrities, are intricately woven into the latent space of the model, making them difficult to remove. Unlike simpler concepts such as nudity, which can often be managed by filtering outputs, identities are deeply embedded in the model’s knowledge, requiring more sophisticated techniques for effective erasure.

Motivation. A widely used unlearning approach in diffusion models involves fine-tuning the original model by conditioning the noise estimate on the target concept to be removed, guiding it to align with the unconditional estimate Kumari et al. (2023); Gandikota et al. (2023); Kim et al. (2023). All the existing methods use the L2 loss function during fine-tuning. The choice of loss function is critical for guiding the model’s ability to selectively remove undesirable concepts, such as nudity

054 or artistic styles, while preserving image quality in retain set data. We empirically observe that L2
 055 loss is limited in handling complex, multi-modal distributions within the model’s latent space. This
 056 motivated us to explore other loss functions such as Jensen-Shannon (JS) divergence and Cramer-
 057 Wold distance. Both provide a nuanced way to align the model’s output distribution with the target
 058 distribution, allowing for more effective unlearning without catastrophic forgetting.

059 **Our Contributions.** In this work, we empirically demonstrate that employing L2 loss as the default
 060 loss function in the unlearning (fine-tuning) setup results in suboptimal performance in current un-
 061 learning methods. We study the impact of employing advanced loss functions such as JS divergence
 062 and Cramer-Wold distance within a knowledge distillation framework. We hypothesize that L2 loss
 063 provides a simple pointwise correction, whereas, JS divergence and Cramer-Wold distance enable
 064 a more robust alignment of distributions within the latent space of diffusion models. These loss
 065 functions are better at handling complex, multi-modal latent spaces. This allows for the removal of
 066 targeted concepts while retaining the model’s generative capabilities. By leveraging these loss func-
 067 tions, we ensure that unlearning is adaptable across different scenarios and minimizes interference
 068 with non-targeted concepts, resulting in a more versatile and reliable unlearning process.

069 We demonstrate how JS divergence and Cramer-Wold distance can be effectively integrated into
 070 the unlearning process to guide the removal of targeted concepts while preserving image quality
 071 in the remaining dataset. We provide comprehensive experiments showing that our approach not
 072 only outperforms existing L2-based methods in concept removal but also mitigates the risk of catas-
 073 trophic forgetting, thereby maintaining the model’s versatility across various unlearning scenarios.
 074 By introducing these loss functions in the knowledge distillation framework and demonstrating their
 075 effectiveness to erase concepts, we offer a significant step forward in the development of safe and
 076 reliable unlearning techniques for diffusion models.

077 2 RELATED WORK

078 We discuss some of the major challenges in ensuring effective moderation of content generated by
 079 text-to-image diffusion models. These challenges include handling NSFW and restricted content,
 080 addressing rights and usage concerns. We review how prior research has approached these issues.

081 **NSFW and Restricted Content.** Diffusion models can be misused to generate inappropriate im-
 082 ages, including violent or explicit content. Some methods attempt to mitigate this by filtering train-
 083 ing data or employing post-processing safety checks during inference Gandhi et al. (2020); Nichol
 084 et al. (2022); Bedapudi; Rando et al.; Schramowski et al. (2023). While effective, these methods can
 085 be bypassed when open-source code and model weights are publicly accessible Sharma et al. (2024).
 086 A more difficult-to-circumvent approaches involve altering the model’s knowledge by modifying at-
 087 tention weights Gandikota et al. (2023), through fine-tuning Zhang et al. (2024b), or with continual
 088 learning Heng & Soh (2024). However, they have been shown to be prone to partial diffusion based
 089 attacks in Sharma et al. (2024). Our method proposes a more effective solution to erase unwanted
 090 concepts that also performs well in the recently introduced unlearning metrics Sharma et al. (2024).

091 **Rights and Usage Concerns.** Stable diffusion models trained on extensive datasets such as LAION-
 092 5B Schuhmann et al. (2022) have been found to infringe on copyrights related to artistic styles,
 093 largely because of their tendency to memorize copyrighted material in the training data Zhang et al.
 094 (2024c). Under U.S. copyright law ¹, the concept of *Fair Use* allows others to use copyrighted ma-
 095 terial only if it is transformed in a way that differs significantly from the original. Legal precedents
 096 have shown that structural similarity to the original work can result in infringement claims.

097 Somepalli et al. (2023); Wang et al. (2024) examined how diffusion models can infringe on copy-
 098 rights through memorization, enabling the generation of existing posters, artworks, and other images
 099 that may be protected by copyright. To address this, some studies have suggested adding minor per-
 100 turbations, acting as watermarks, into images to prevent their memorization by diffusion models Cui
 101 et al. (2023); Zhao et al. (2023). However, these watermarks can be readily eliminated using tech-
 102 niques like denoising or blurring. Another approach involves post-training model adjustments to
 103 delete undesired concepts, such as modifying model weights to remove specific styles Gandikota
 104 et al. (2023); Heng & Soh (2024); Kim et al. (2023); Kumari et al. (2023); Zhang et al. (2024b).

105 ¹<https://www.copyright.gov/title17/>

Unlearning in Diffusion Models. Unlearning in diffusion models presents a challenge due to the intricate and interconnected nature of their latent space representations. Concepts within these models are not independent but are woven into a larger knowledge distribution. Consequently, the removal of specific elements, like copyrighted logos or identifiable human faces, without compromising the model’s ability to produce high-quality images is difficult. Current leading methods in the field involve fine-tuning the diffusion model by adjusting attention heads or using distillation techniques that condition the noise estimate on the concept to be removed Gandikota et al. (2023); Heng & Soh (2024); Kim et al. (2023); Kumari et al. (2023); Zhang et al. (2024b). However, recent findings by Sharma et al. (2024) indicate that many of these approaches Kumari et al. (2023); Gandikota et al. (2023) do not completely erase the concepts, as they can still be generated using partial diffusion probing.

3 PRELIMINARIES

3.1 UNLEARNING IN DIFFUSION MODELS

Unlearning aims to remove specific learned concepts from diffusion models while preserving the model’s overall performance. In diffusion models, the goal is to eliminate targeted knowledge, such as biases or harmful content, without degrading the image generation quality. The unlearning process can be expressed through the following general equation:

$$\theta_{\text{unlearned}} = \theta_{\text{fully trained}} - \eta \sum_{t=1}^T \nabla_{\theta} L_t(\theta; \mathcal{D}_r, x_t, c), \quad (1)$$

where $\theta_{\text{unlearned}}$ is the model’s parameter after the unlearning process, $\theta_{\text{fully trained}}$ is the model’s parameter before unlearning, η is the learning rate, T is the total number of diffusion steps, c denotes the conditioning using a particular modality, and $L_t(\theta; \mathcal{D}_r, x_t, c)$ is the loss function at step t .

3.2 UNLEARNING VIA DISTILLATION

Several existing methods have adopted different forms of knowledge distillation for unlearning in classification models, regression models, language models and diffusion models. We define the general form of the unlearning objective as:

$$\mathcal{L} = \mathcal{D}(\epsilon_{\theta_{\text{student}}}(z_t, c_s, t), \text{sg}(\epsilon_{\theta_{\text{student}}}(z_t, t))), \quad (2)$$

where $\mathcal{D}(\cdot, \cdot)$ is a generic distance function or loss metric. This could be the L^2 norm, L^1 norm, Kullback-Leibler (KL) divergence, or any other suitable loss function. Here, $\epsilon_{\theta_{\text{student}}}(z_t, c_s, t)$ is the noise estimate conditioned on the target concept c_s , and $\epsilon_{\theta_{\text{student}}}(z_t, t)$ is the unconditional noise estimate. The stop-gradient operation, $\text{sg}(\cdot)$, is applied to prevent gradient computation for the unconditional noise estimate, ensuring that updates are made solely based on the conditioned estimate.

The teacher model is then updated using the student’s parameters to smooth out abrupt changes:

$$\theta_{\text{teacher}} \leftarrow \phi(\theta_{\text{teacher}}, \theta_{\text{student}}), \quad (3)$$

where $\phi(\cdot, \cdot)$ denotes the update rule that integrates the student’s potentially drastic latent updates into a more stable adjustment for the teacher model. This update process reflects the student’s learning while ensuring that the changes to the teacher model remain gradual and controlled, avoiding harsh updates that may destabilize the model.

4 PROPOSED WORK

We formally define the key mathematical properties for three distances Cramér-Wold Cramér & Wold (1936) distance, Jensen-Shannon Divergence Lin (1991); Rao & Nayak (1985), and L2 distance experimented in this work. To our knowledge this is the first work to mathematically and empirically analyze the significance of each distances in the context of unlearning in diffusion models. We introduce four lemmas by analyzing several mathematical principles such as functional

analysis (norms and decomposition), measure theory (distributional differences), information theory (entropy and divergence measures), and orthogonal projections. We will use these principles to rigorously demonstrate the limitations of L2 distance and the effectiveness of JS divergence and Cramér-Wold distance in capturing concept-specific changes in a diffusion model.

4.1 LIMITATIONS OF L2 DISTANCE IN UNLEARNING

We identify two main limitations of using L2 distance in unlearning. First, we demonstrate that L2 distance is highly sensitive to irrelevant dimensions. To this end we formally show:

Lemma 1. (L2 Distance is Sensitive to Irrelevant Dimensions)

Let $\mathbf{z}_T, \mathbf{z}_S \in \mathbb{R}^n$ be two n -dimensional vectors representing the teacher and student embeddings respectively, and let the concept \mathbf{c} to be unlearned be represented in a subset of dimensions $\mathcal{D}_c \subseteq \{1, 2, \dots, n\}$. The L2 distance between \mathbf{z}_T and \mathbf{z}_S is non-zero even when the concept is fully unlearned, as long as \mathbf{z}_T and \mathbf{z}_S differ in irrelevant dimensions, $\mathcal{D}_{nc} = \{1, 2, \dots, n\} \setminus \mathcal{D}_c$.

The proof for Lemma 1 shows that the L2 distance between vectors \mathbf{z}_T and \mathbf{z}_S can be decomposed into concept-related and non-concept-related parts. If the concept is fully unlearned, the contribution from concept-related components becomes zero. However, if differences exist in non-concept-related components, the overall L2 distance remains positive, indicating that the L2 distance is not zero even when the concept is unlearned. The detailed proof is added in the Appendix.

Next we show that L2 distance have difficulty capturing correlated conceptual changes, which is crucial property for any unlearned system.

Lemma 2. (L2 Distance Fails to Capture Correlated Conceptual Changes)

Let $\mathbf{z}_T, \mathbf{z}_S \in \mathbb{R}^n$ be two n -dimensional vectors. Suppose the concept \mathbf{c} is represented as a linear combination of multiple correlated dimensions in \mathbf{z}_T . If the unlearning process changes these dimensions uniformly in \mathbf{z}_S , then the L2 distance will overestimate or underestimate the actual conceptual difference.

The Lemma 2 argues that if a concept is represented by multiple correlated dimensions in \mathbf{z}_T , and these dimensions are uniformly altered during the unlearning process in \mathbf{z}_S , the L2 distance will not accurately reflect the true conceptual change. This is because L2 distance treats each dimension independently, failing to capture the correlation between dimensions. As a result, the L2 distance may either overestimate or underestimate the difference depending on the nature and scale of the uniform changes, the detailed proof can be found in appendix.

4.2 QUANTIFYING CONCEPTUAL DIFFERENCES WITH JS DIVERGENCE

We show JS divergence is well equipped to quantify conceptual differences.

Lemma 3. (Jensen-Shannon Divergence Accurately Quantifies Conceptual Differences)

Let \mathbf{z}_T and \mathbf{z}_S be n -dimensional embeddings represented as probability distributions P and Q over the same n -dimensional space. Suppose the concept \mathbf{c} is encoded in a subset of dimensions $\mathcal{D}_c \subseteq \{1, 2, \dots, n\}$. If P and Q differ only in the concept dimensions \mathcal{D}_c and are identical in the remaining dimensions $\mathcal{D}_{nc} = \{1, 2, \dots, n\} \setminus \mathcal{D}_c$, the Jensen-Shannon (JS) divergence will accurately capture the concept difference while remaining invariant to changes in \mathcal{D}_{nc} .

The proof for Lemma 3 demonstrates that the Jensen-Shannon (JS) divergence is a better measure for concept unlearning compared to L2 distance. It starts by defining the JS divergence between two probability distributions P and Q , and decomposes these distributions into concept-related and non-concept-related components. By showing that if the non-concept-related parts of P and Q are equal, the JS divergence becomes zero for these dimensions, it highlights that JS divergence is invariant to changes in irrelevant dimensions. In contrast, L2 distance would still capture these differences, making it sensitive to irrelevant variations. Therefore, JS divergence accurately captures conceptual changes, while L2 distance may overestimate the differences. The detailed proof is provided in the Appendix.

Table 1: We compare the CLIP scores obtained by the CW, JS and L2 unlearning methods. We show the mean CLIP score by unlearning the following concepts: *baby*, *narendra modi*, *elon musk*, *amitabh bachchan*, *nike*, *nudity*, *pablo picasso*, *vincent van gogh*. *lower is better*↓.

| Steps | SD 1.4 | CW | JS | L2 |
|-------|---------|----------------|----------------|----------------|
| 1500 | 29.7618 | 23.4975 | 22.9702 | 23.4436 |
| 1400 | 29.7618 | 23.4934 | 22.6502 | 23.0212 |
| 1300 | 29.7618 | 24.1404 | 24.4098 | 23.4957 |
| 1200 | 29.7618 | 25.2973 | 24.8594 | 23.0008 |
| 1100 | 29.7618 | 25.5798 | 25.8554 | 26.3378 |

Table 2: Comparing the unlearning performance in terms of *CCS*, and *CRS* score. We compare the CW and JS with L2 loss based method. We evaluate the effectiveness of concept erasure for three types of concepts: *celebrity*, *baby*, and *artistic style*.

| Concept Erased | Prompt | <i>CCS</i> | | | <i>CRS</i> | | |
|------------------|--------------------------------|------------|------|------|------------|-------|-------|
| | | CW | JS | L2 | CW | JS | L2 |
| amitabh bachchan | amitabh bachchan | 0.74 | 0.61 | 0.74 | 0.05 | 0.05 | 0.04 |
| baby | baby with teddy bear | 0.62 | 0.62 | 0.63 | 0.02 | 0.02 | 0.03 |
| vincent van gogh | sunflowers by vincent van gogh | 0.51 | 0.42 | 0.52 | 0.009 | 0.008 | 0.008 |

4.3 CAPTURING HIGH-ORDER CONCEPTUAL CORRELATIONS WITH CRAMÉR-WOLD DISTANCE

Cramér-Wold distance has ability to capture higher-order correlations and joint distributional changes between dimensions, which can be formally represented as follows:

Lemma 4. (Cramér-Wold Distance Captures High-Order Conceptual Correlations)

Let $\mathbf{z}_T, \mathbf{z}_S \in \mathbb{R}^n$ be two vectors, and let P, Q be the corresponding distributions over these vectors. If the concept \mathbf{c} is represented by correlations between multiple dimensions, then the Cramér-Wold distance between the distributions of \mathbf{z}_T and \mathbf{z}_S will be zero if and only if \mathbf{c} has been fully unlearned, even if the individual L2 distances in some dimensions are non-zero.

The proof for Lemma 4 establishes that Cramér-Wold (CW) distance is a more suitable measure for detecting conceptual changes compared to L2 distance. It begins by defining the CW distance as the supremum over 1-dimensional linear projections of the distributions P and Q . It then shows that CW distance captures correlations by analyzing projections aligned with the concept subspace. If a concept is represented by correlated dimensions, changes in these dimensions will be detected by CW distance but may be missed by L2 distance, which only considers dimension-wise differences. Finally, the proof highlights that CW distance is invariant to changes in irrelevant projections, making it robust in measuring concept changes across different subspaces. The detailed proof is provided in the Appendix.

5 EXPERIMENTS AND RESULTS

We show that both Cramér-Wold distance and Jensen-Shannon (JS) divergence outperform L2 distance in a knowledge distillation setup for unlearning in diffusion models. Our approach is compared against existing state-of-the-art methods ESD-u, ESD-x Gandikota et al. (2023), Ablating Concepts Kumari et al. (2023), and SDD Kim et al. (2023), which also utilize a knowledge distillation framework but rely on L2 distance as their default loss function. We choose SDD Kim et al. (2023) to represent all four distillation-based unlearning methods, as it has already been demonstrated to outperform ESD-u, ESD-x, and Ablating concepts in Sharma et al. (2024). The L2 loss results are basically the results in SDD. We denote Cramér-Wold results as **CW**, JS Divergence results as **JS**, and L2 loss-based method as **L2** in our discussions.

Table 3: CLIP directional similarity score based comparison of the proposed CW and JS methods with L2 loss based method. Both CW and JS perform better erasure as compared to the L2 based unlearning. *higher is better*↑.

| Concept Erased | Prompt | CW | JS | L2 |
|------------------|-------------------------------------|---------------|---------------|---------------|
| baby | a baby playing with teddy bear | 0.1289 | 0.1874 | 0.1363 |
| narendra modi | narendra modi in a park | 0.1375 | 0.1876 | 0.0996 |
| elon musk | elon musk shaking hand with a girl | 0.2085 | 0.1802 | 0.1348 |
| amitabh bachchan | amitabh bachchan in blue background | 0.1682 | 0.2138 | 0.1738 |
| nike | a shoe with nike logo | 0.0926 | 0.0585 | 0.1058 |
| mean score | | 0.1328 | 0.1463 | 0.1284 |

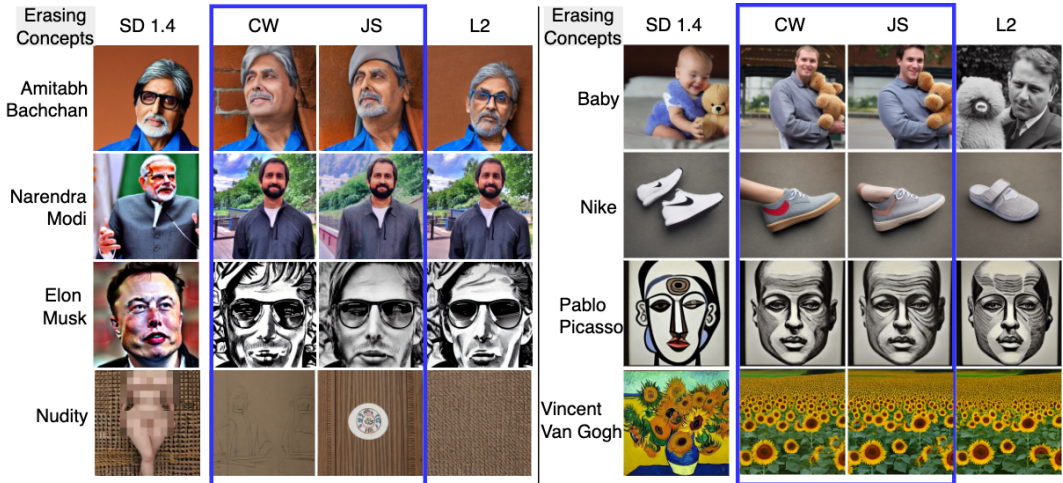


Figure 1: The concept erasure results are shown for the CW, JS, and L2 methods for a variety of erasure tasks. We observe similar results in most of the categories except in *Nike* concept erasure, where CW and JS generated images of shoe without *Nike* logo, whereas, L2 generated footwear that is not a shoe but looks like *crocs*.

Evaluation of Concept Erasure. We assess the effectiveness of unlearning using four metrics: the CLIP score, CLIP directional similarity score Kumari et al. (2023), Concept Confidence Score (CCS), and Concept Retrieval Score (CRS) Sharma et al. (2024). Additionally, we provide qualitative results for both the erased and retained concepts, allowing for a fair comparison of the various unlearning methods.

Experiment Setting. We assess the performance of concept erasure across the following categories: art style, logo, identity, and NSFW content, utilizing Stable Diffusion 1.4 (SD 1.4). The experiments were carried out with NVIDIA A6000 48GB GPU. We observe that different methods give optimal results at varying iterations. Therefore, we present the results at multiple iteration intervals to ensure a fair comparison.

Quantitative Analysis. We show the mean CLIP scores for the 8 *erased concepts* from the SD 1.4 model in Table 1. The CLIP score is a metric used to measure how well the generated images semantically match their target concepts. For the *erased concepts*, lower CLIP score indicates a reduced similarity between the generated images and the given prompt, suggesting more effective concept erasure. For example, after erasing the concept “Amitabh Bachchan”, the model would not generate images resembling to *Amitabh Bachchan*. In Table 1, JS method (ours) consistently outperforms CW (ours) and L2 (existing method) at steps 1500, 1400, and 1300. At step 1200, L2 is better, while at step 1100, CW gives better score. JS method exhibits particularly strong performance, achieving the lowest CLIP scores in the majority of the steps. Overall, we observe CW and JS either maintains similar performance or improves upon the L2 method.

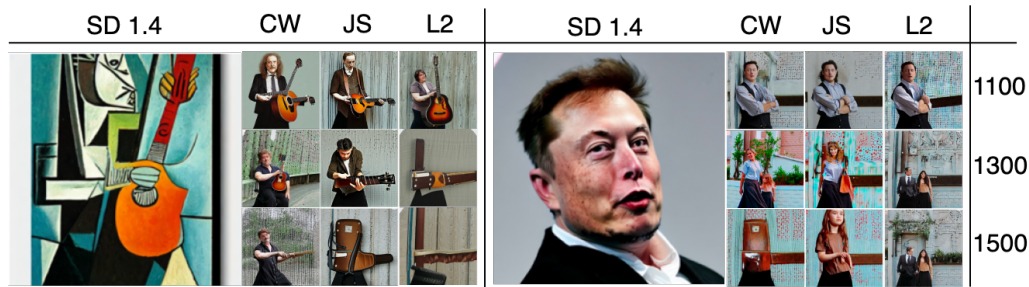


Figure 2: Comparing the unlearning results of CW, JS, and L2 after erasing the concepts of *Pablo Picasso* and *Elon Musk* from the SD 1.4. Prompts: *the old guitarist by Pablo Picasso*, *Elon Musk*. We observe that at 1100 step, both CW and JS generate old guitarist but without the style of Pablo Picasso but L2 generates a visibly young guitarist.



Figure 3: We show the progression of concept unlearning across iterations (100 to 1500) for CW, JS, and L2 methods.

We examine the CLIP directional similarity scores across various domains, with the results presented in Table 3. Five different concepts were selected, and we applied unlearning to the SD 1.4 model for each. The CLIP directional similarity scores for JS, CW, and L2 were then evaluated, where a higher score indicates more effective unlearning. For the concept of “baby,” JS achieved the highest score of 0.1874. Similarly, CW obtained a score of 0.2085 when erasing the concept of “Elon Musk,” outperforming both L2 and JS. On average, JS outperforms both CW and L2. Overall, the proposed JS and CW methods perform better than L2 in terms of the CLIP directional similarity score.

We evaluate the methods CW, JS and L2 based on the adversarial recovery attacks Concept Confidence Score (CCS), Concept Retrieval Score (CRS). We conduct the erasure for three distinct categories: *celebrity* (*Amitabh Bachchan*), *object* (*baby*), and *artistic style* (*vincent van gogh*). We observe similar performance for all three methods CW, JS, and L2 in Table 2. This indicates that the change in loss function doesn’t influence the ability of the unlearned model against adversarial recovery attacks such as CCS and CRS.

Qualitative Analysis. We show the visual results for erasing varieties of concepts: *artistic style*, *celebrity*, *baby/children*, *nudity*, *logo*, and *harmful contents*. In Figure 1, we show the unlearning results of CW, JS, and L2. In most of the concept erasure requests, we observe similar results except for erasure of *Nike* logo. In *Nike* concept erasure, CW and JS generate images of shoe with a different type of logo, whereas, L2 generates a footwear that is not a shoe but looks like cros. We show another comparison between these methods in Figure 2. We erase the concepts of *Pablo Picasso* and *Elon Musk* from SD 1.4. We observe CW and JS generate old guitarist without the style of Pablo Picasso but L2 generates a young guitarist.

Unlearning at Different Iterations. Figure 3 illustrates the unlearning process across different iterations in JS, CW, and L2. The visuals provides insights into the dynamics of concept erasure when using different loss functions like JS, CW, and L2. After 800 steps, CW(ours) and JS(ours) generate different looking image than *Amitabh Bachchan* while L2 still generates an image similar to *Amitabh Bachchan*. In the final step (step 1500), L2 fails to generate human face while CW and JS generate human face that is not similar to *Amitabh Bachchan*.

Visualization at the Teacher and Student Model. Figure 4 compares the outputs of the teacher and student models during the unlearning process in JS, CW, and L2. The comparison highlights the

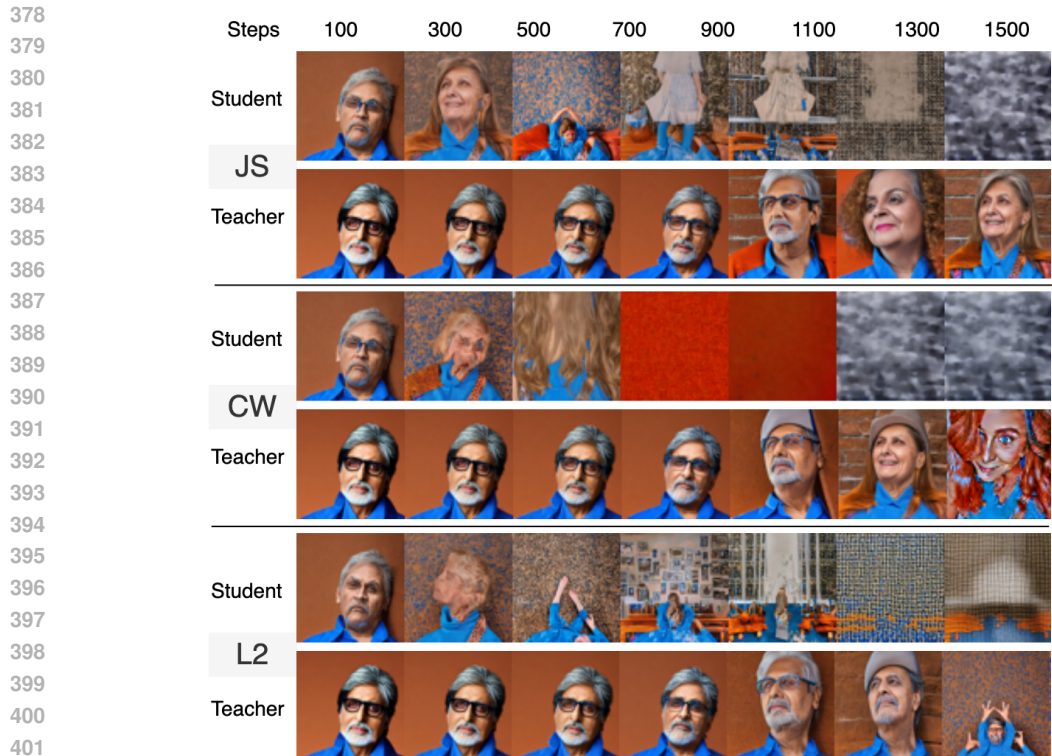


Figure 4: We compare the change in generated outputs by the teacher and student models in CW, JS, and L2. We show the outputs of teacher and student models at different iterations (100 to 1500) to assess the impact of using different loss functions.

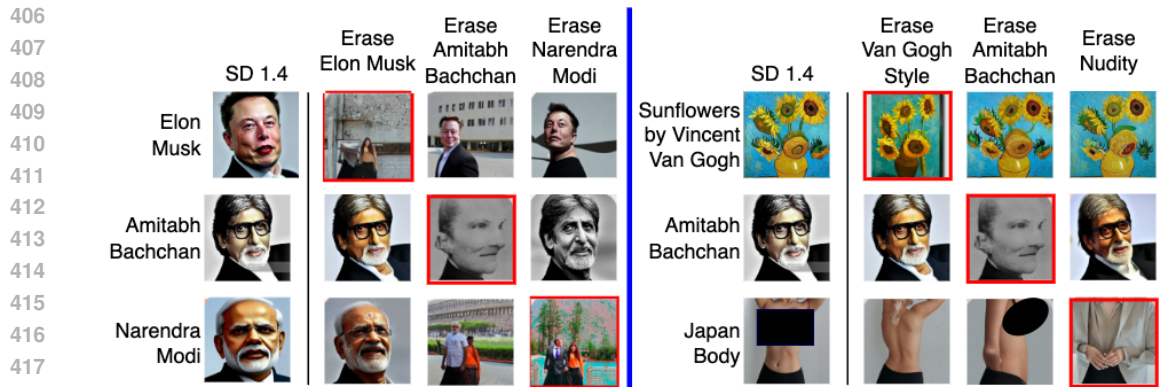


Figure 5: We show the results produced by the proposed method (JS) over the unlearned and retention set of concepts. The diagonal images are the generated for the *erasure concepts* and the remaining images are generated for the *retained concepts* (for the corresponding unlearned model). It shows our method preserves the general capabilities of the diffusion model after unlearning.

effectiveness of teacher-student framework in transferring the unlearning process while maintaining overall image quality. The student that is given an empty prompt is naturally producing noisy image for all the three methods. The teacher model which is eventually used for unlearning shows that JS and CW learn finer variations for the erasing concepts as compared to the L2 method.

Observing the Results for the Retained Concepts. To assess whether the unlearning methods maintain the model’s overall capabilities post-unlearning, we present results on the retained concepts for the modified models. Figure 5 displays a diverse set of generated images unrelated to the erased concepts. The images along the diagonal represent the concepts meant to be removed, while the

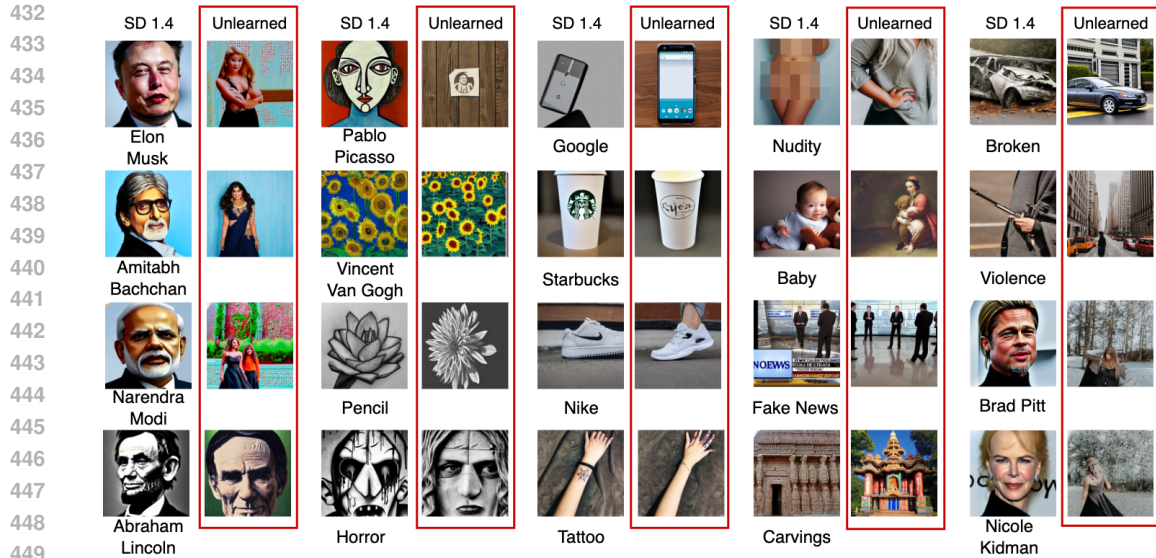


Figure 6: Before and after unlearning results for a variety of unlearning requests that cover the usecases for NSFW, restricted content, rights, and usage concerns. The results are shown for JS method.

454 remaining images correspond to the concepts the model should retain. It is evident that the proposed
455 JS method effectively preserves the retained concepts while successfully forgetting the targeted ones.

456 **Concepts Generated Before and After Unlearning.** Figure 6 illustrates the results of our unlearn-
457 ing process across various concepts, both before and after applying the method. The comparison
458 highlights the model’s capability to effectively remove specific concepts while preserving overall
459 image quality and coherence. The examples include facial unlearning (e.g., Elon Musk, Amitabh
460 Bachchan, Narendra Modi, Brad Pitt, Nicole Kidman), artistic styles (e.g., Pablo Picasso, Vincent
461 van Gogh, pencil art, tattoos), logos (e.g., Google, Starbucks, Nike), and other miscellaneous cate-
462 gories (e.g., nudity, children, fake news, temple carvings, violence, broken car). Additional results
463 are shown in Figure 7, Figure 8, and Figure 9.

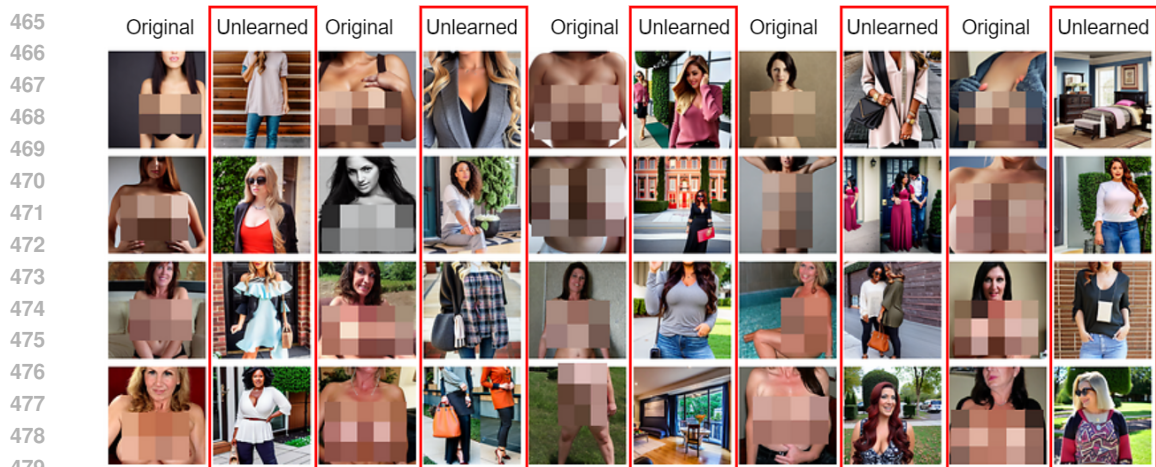
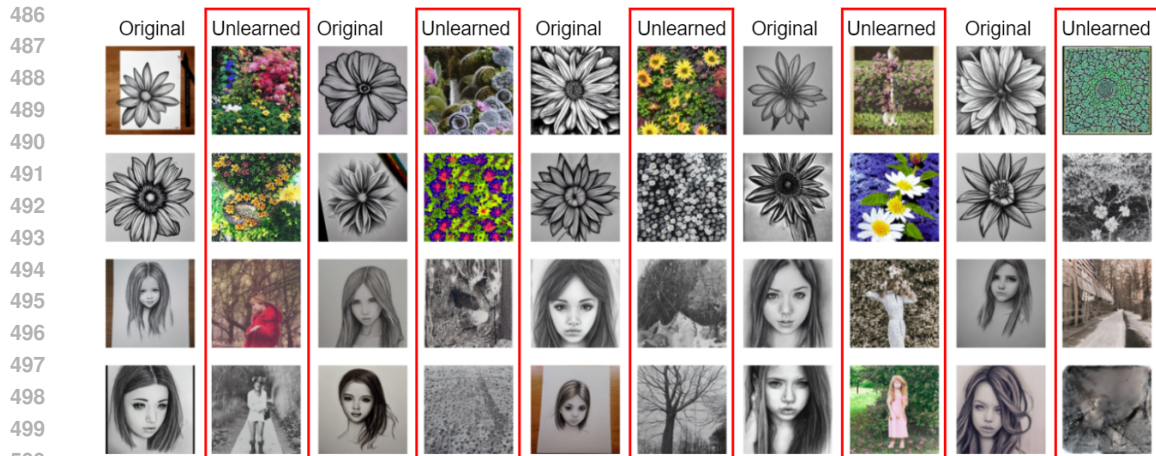


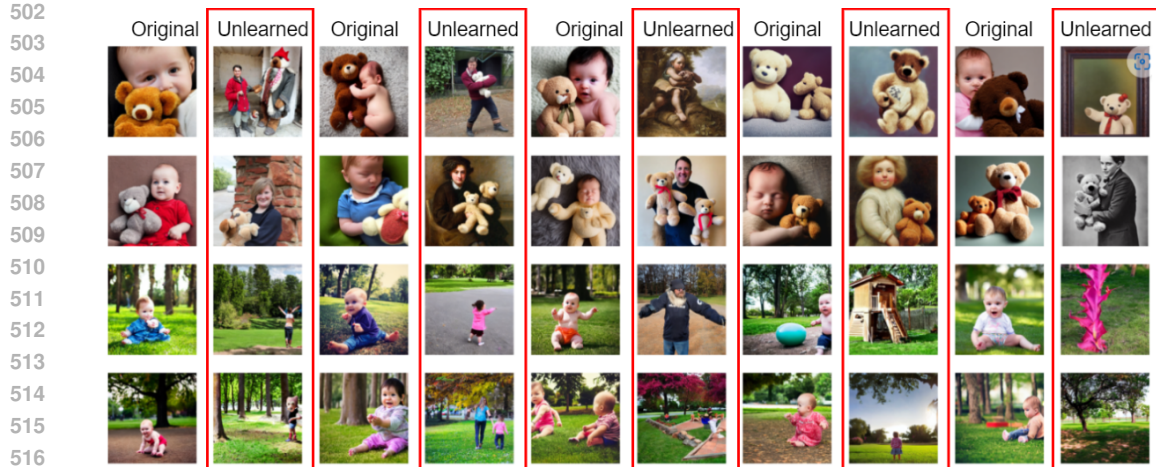
Figure 7: Before and after unlearning results while erasing the concept of *Nudity* (JS method).

482 6 CONCLUSION

483 In this paper, we proposed a set of robust methodologies for concept unlearning in diffusion models,
484 utilizing advanced loss functions like Cramer-Wold distance and Jensen-Shannon (JS) divergence
485



500 Figure 8: Before and after unlearning results while erasing the concept of *Pencil Art* (JS method).



517 Figure 9: Before and after unlearning results while erasing the concept of *Baby* (JS method).

519 within a knowledge distillation framework. Our approach demonstrated improved performance over
520 L2 loss based method over diverse concept removal categories. Extensive experiments and math-
521 ematical analysis confirmed the versatility and effectiveness of our unlearning method, paving the
522 way for more controlled and responsible applications of text-to-image diffusion models.

524 REFERENCES

- 526 Praneeth Bedapudi. Nudenet: Neural nets for nudity detection and censoring, 2022. URL
527 <https://github.com/notAI-tech/NudeNet>.
- 528 Harald Cramér and Herman Wold. Some theorems on distribution functions. *Journal of the London*
529 *Mathematical Society*, 1(4):290–294, 1936.
- 531 Yingqian Cui, Jie Ren, Han Xu, Pengfei He, Hui Liu, Lichao Sun, Yue Xing, and Jiliang Tang.
532 Diffusionshield: A watermark for copyright protection against generative diffusion models. *arXiv*
533 *preprint arXiv:2306.04642*, 2023.
- 534 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
535 *in neural information processing systems*, 34:8780–8794, 2021.
- 537 Shreyansh Gandhi, Samrat Kokkula, Abon Chaudhuri, Alessandro Magnani, Theban Stanley, Be-
538 hzad Ahmadi, Venkatesh Kandaswamy, Omer Ovenc, and Shie Mannor. Scalable detection of
539 offensive and non-compliant content/logo in product images. In *Proceedings of the IEEE/CVF*
winter conference on applications of computer vision, pp. 2247–2256, 2020.

- 540 Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts
541 from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer
542 Vision*, pp. 2426–2436, 2023.
- 543 Aditya Golatkar, Alessandro Achille, Luca Zancato, Yu-Xiang Wang, Ashwin Swaminathan, and
544 Stefano Soatto. Cpr: Retrieval augmented generation for copyright protection. In *Proceedings of
545 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12374–12384, 2024.
- 546 Jiatao Gu, Shuangfei Zhai, Yizhe Zhang, Joshua M Susskind, and Navdeep Jaitly. Matryoshka
547 diffusion models. In *The Twelfth International Conference on Learning Representations*, 2023.
- 548 Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep
549 generative models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 550 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in
551 neural information processing systems*, 33:6840–6851, 2020.
- 552 Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration
553 models. *Advances in Neural Information Processing Systems*, 35:23593–23606, 2022.
- 554 Sanghyun Kim, Seohyeon Jung, Balhae Kim, Moonseok Choi, Jinwoo Shin, and Juho Lee. Towards
555 safe self-distillation of internet-scale text-to-image diffusion models, 2023.
- 556 Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan
557 Zhu. Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF
558 International Conference on Computer Vision*, pp. 22691–22702, 2023.
- 559 Jianhua Lin. Divergence measures based on the shannon entropy. *IEEE Transactions on Information
560 theory*, 37(1):145–151, 1991.
- 561 Shilin Lu, Zilan Wang, Leyang Li, Yanzhu Liu, and Adams Wai-Kin Kong. Mace: Mass concept
562 erasure in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision
563 and Pattern Recognition*, pp. 6430–6440, 2024.
- 564 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob
565 Mcgrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and
566 editing with text-guided diffusion models. In *International Conference on Machine Learning*, pp.
567 16784–16804. PMLR, 2022.
- 568 Javier Rando, Daniel Paleka, David Lindner, Lennart Heim, and Florian Tramer. Red-teaming the
569 stable diffusion safety filter. In *NeurIPS ML Safety Workshop*.
- 570 C Rao and T Nayak. Cross entropy, dissimilarity measures, and characterizations of quadratic
571 entropy. *IEEE Transactions on Information Theory*, 31(5):589–593, 1985.
- 572 Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion:
573 Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF
574 Conference on Computer Vision and Pattern Recognition*, pp. 22522–22531, 2023.
- 575 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
576 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
577 open large-scale dataset for training next generation image-text models. *Advances in Neural
578 Information Processing Systems*, 35:25278–25294, 2022.
- 579 Aakash Sen Sharma, Niladri Sarkar, Vikram Chundawat, Ankur A Mali, and Murari Mandal. Un-
580 learning or concealment? a critical analysis and evaluation metrics for unlearning in diffusion
581 models. *arXiv preprint arXiv:2409.05668*, 2024.
- 582 Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion
583 art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the
584 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6048–6058, 2023.
- 585 Haonan Wang, Qianli Shen, Yao Tong, Yang Zhang, and Kenji Kawaguchi. The stronger the diffu-
586 sion model, the easier the backdoor: Data poisoning to induce copyright breaches without adjust-
587 ing finetuning pipeline. *arXiv preprint arXiv:2401.04136*, 2024.
- 588
589
590
591
592
593

594 Gong Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not: Learn-
595 ing to forget in text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on*
596 *Computer Vision and Pattern Recognition*, pp. 1755–1764, 2024a.

597
598 Gong Zhang, Kai Wang, Xingqian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not: Learn-
599 ing to forget in text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on*
600 *Computer Vision and Pattern Recognition*, pp. 1755–1764, 2024b.

601 Yang Zhang, Teoh Tze Tzun, Lim Wei Hern, and Kenji Kawaguchi. On copyright risks of text-to-
602 image diffusion models. In *ECCV 2024 Workshop The Dark Side of Generative AIs and Beyond*,
603 2024c.

604
605 Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Ngai-Man Cheung, and Min Lin. A recipe for
606 watermarking diffusion models. *arXiv preprint arXiv:2303.10137*, 2023.

607 608 A APPENDIX

609 610 A.1 PROOF OF LEMMA1

611
612 *Proof.* The L2 distance between $\mathbf{z}_T = [z_{T,1}, z_{T,2}, \dots, z_{T,n}]$ and $\mathbf{z}_S = [z_{S,1}, z_{S,2}, \dots, z_{S,n}]$ is given
613 by:

$$614 \quad 615 \quad 616 \quad 617 \quad 618 \quad \text{L2}(\mathbf{z}_T, \mathbf{z}_S) = \sqrt{\sum_{i=1}^n (z_{T,i} - z_{S,i})^2} \quad (4)$$

619 We can partition the sum into concept-related and non-concept-related dimensions:

$$620 \quad 621 \quad 622 \quad 623 \quad 624 \quad \text{L2}(\mathbf{z}_T, \mathbf{z}_S) = \sqrt{\sum_{i \in \mathcal{D}_c} (z_{T,i} - z_{S,i})^2 + \sum_{i \in \mathcal{D}_{nc}} (z_{T,i} - z_{S,i})^2} \quad (5)$$

$$625 \quad 626 \quad 627 \quad 628 \quad 629 \quad (6)$$

630 If the concept c is unlearned, then $z_{T,i} = z_{S,i}$ for all $i \in \mathcal{D}_c$. Hence:

$$631 \quad 632 \quad 633 \quad \sum_{i \in \mathcal{D}_c} (z_{T,i} - z_{S,i})^2 = 0. \quad (7)$$

634 However, if $z_{T,i} \neq z_{S,i}$ for any $i \in \mathcal{D}_{nc}$, then:

$$635 \quad 636 \quad 637 \quad \sum_{i \in \mathcal{D}_{nc}} (z_{T,i} - z_{S,i})^2 > 0, \quad (8)$$

638 implying:

$$639 \quad 640 \quad \text{L2}(\mathbf{z}_T, \mathbf{z}_S) > 0 \quad (9)$$

641 Thus, L2 distance is not zero even if the concept is fully unlearned, as long as there are differences
642 in irrelevant dimensions. \square

643 644 A.2 PROOF OF LEMMA2

645
646 *Proof.* Let the concept c be represented as a linear combination of dimensions i, j, k such that:

$$647 \quad \mathbf{c} = \alpha z_{T,i} + \beta z_{T,j} + \gamma z_{T,k}. \quad (10)$$

Suppose unlearning results in a proportional decrease in these dimensions in \mathbf{z}_S , i.e.,

$$z_{S,i} = z_{T,i} - \delta, \quad z_{S,j} = z_{T,j} - \delta, \quad z_{S,k} = z_{T,k} - \delta, \quad (11)$$

for some constant δ . The L2 distance is:

$$\text{L2}(\mathbf{z}_T, \mathbf{z}_S) = \sqrt{(z_{T,i} - (z_{T,i} - \delta))^2 + (z_{T,j} - (z_{T,j} - \delta))^2 + (z_{T,k} - (z_{T,k} - \delta))^2} = \delta\sqrt{3} \quad (12)$$

However, the true change in the concept is along \mathbf{c} , and the L2 distance does not reflect the conceptual change accurately unless $\alpha = \beta = \gamma$. Therefore, L2 distance fails to correctly capture the change in correlated dimensions. \square

A.3 PROOF OF LEMMA3

Proof. Let's recall the definition of JS Divergence 1. *Definition of JS Divergence:* Let $P = [p_1, p_2, \dots, p_n]$ and $Q = [q_1, q_2, \dots, q_n]$ be the probability distributions derived from the teacher and student representations such that:

$$\sum_{i=1}^n p_i = 1, \quad \sum_{i=1}^n q_i = 1, \quad p_i, q_i \geq 0 \quad (13)$$

The Jensen-Shannon (JS) divergence between P and Q is defined as:

$$\text{JS}(P, Q) = \frac{1}{2}KL(P\|M) + \frac{1}{2}KL(Q\|M), \quad (14)$$

where $M = \frac{P+Q}{2}$ and KL is the Kullback-Leibler divergence:

$$KL(P\|M) = \sum_{i=1}^n p_i \log \frac{p_i}{m_i}, \quad KL(Q\|M) = \sum_{i=1}^n q_i \log \frac{q_i}{m_i}, \quad m_i = \frac{p_i + q_i}{2} \quad (15)$$

2. *Partition into Concept and Non-Concept Dimensions:* Let \mathcal{D}_c denote the set of concept-related dimensions and \mathcal{D}_{nc} denote the irrelevant dimensions. Decompose P and Q as:

$$P = [P_c, P_{nc}], \quad Q = [Q_c, Q_{nc}], \quad (16)$$

where $P_c = \{p_i : i \in \mathcal{D}_c\}$ and $P_{nc} = \{p_i : i \in \mathcal{D}_{nc}\}$ (similarly for Q_c and Q_{nc}). If $P_{nc} = Q_{nc}$, then $p_i = q_i$ for all $i \in \mathcal{D}_{nc}$. For these dimensions, the KL divergence terms are zero:

$$\sum_{i \in \mathcal{D}_{nc}} p_i \log \frac{p_i}{m_i} = 0, \quad \sum_{i \in \mathcal{D}_{nc}} q_i \log \frac{q_i}{m_i} = 0 \quad (17)$$

Thus, the JS divergence simplifies to:

$$\text{JS}(P, Q) = \frac{1}{2} \sum_{i \in \mathcal{D}_c} p_i \log \frac{p_i}{m_i} + \frac{1}{2} \sum_{i \in \mathcal{D}_c} q_i \log \frac{q_i}{m_i} \quad (18)$$

3. *Invariant to Irrelevant Changes:* For $i \in \mathcal{D}_{nc}$, if $p_i = q_i$, the divergence between P and Q in these dimensions will remain zero, regardless of the magnitude of p_i and q_i . In contrast, L2 distance will still consider the differences in \mathcal{D}_{nc} , leading to misleading conclusions about whether the concept has been unlearned.

4. *Effectiveness in Concept Unlearning:* If unlearning is successful and $P_c = Q_c$, then:

$$\text{JS}(P, Q) = 0, \quad (19)$$

702 even if $P_{nc} \neq Q_{nc}$. Thus, JS divergence accurately measures concept alignment by focusing only
 703 on relevant dimensions, whereas L2 distance remains sensitive to changes in irrelevant dimensions.
 704

□

706 A.4 PROOF OF LEMMA4

707
 708 *Proof. 1. Definition of Cramér-Wold Distance:* The Cramér-Wold distance between two distribu-
 709 tions P and Q is defined as the supremum over all 1-dimensional linear projections:
 710

$$711 \quad \text{CW}(P, Q) = \sup_{\|\theta\|=1} \|P_\theta - Q_\theta\|, \quad (20)$$

712 where θ is a unit vector, and P_θ and Q_θ are the 1-dimensional projections of P and Q along θ :
 713

$$714 \quad P_\theta = \theta^T \mathbf{z}_T, \quad Q_\theta = \theta^T \mathbf{z}_S. \quad (21)$$

715
 716 *2. Effectiveness in Capturing Correlations:* If the concept \mathbf{c} is represented by a set of correlated
 717 dimensions \mathcal{D}_c , define a linear combination for the concept:
 718

$$719 \quad \mathbf{c} = \sum_{i \in \mathcal{D}_c} \alpha_i z_{T,i}. \quad (22)$$

720
 721 Consider a projection θ such that it aligns with the concept subspace. If P and Q are identical along
 722 θ , i.e., $P_\theta = Q_\theta$, then:
 723

$$724 \quad \|P_\theta - Q_\theta\| = 0. \quad (23)$$

725
 726 *3. Detecting Higher-Order Correlations:* Unlike L2 distance, which measures dimension-wise dif-
 727 ferences, Cramér-Wold distance takes into account joint distributions and higher-order correlations.
 728 Thus, if unlearning leads to a change in joint correlations but not in individual dimensions, L2 dis-
 729 tance might be zero, while CW distance will detect the conceptual change.
 730

731
 732 *4. Invariant to Irrelevant Projections:* If θ is orthogonal to the concept subspace, then $P_\theta = Q_\theta$
 733 for all such projections, even if $\mathbf{z}_T \neq \mathbf{z}_S$ in irrelevant dimensions. Thus, Cramér-Wold distance
 734 provides a more comprehensive measure of concept unlearning by considering projections along all
 735 directions, ensuring that the concept is completely removed. □
 736

737 A.5 MORE VISUAL RESULTS

738
 739 We show additional visual results of concept erasure in the proposed JS method in Figure 14(eras-
 740 ing *Nike* logo), Figure 15(erasing *Narendra Modi*), Figure 16(erasing *Amitabh Bachchan*), Fig-
 741 ure 17(erasing *Elon Musk*), Figure 18(erasing *child*), Figure 19(erasing *Vincent Van Gogh styled*
 742 *paintings*), Figure 20(erasing *Pablo Picasso styled paintings*).
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809



Figure 10: Unlearning results after erasing the concept of *Brad Pitt* (JS method).

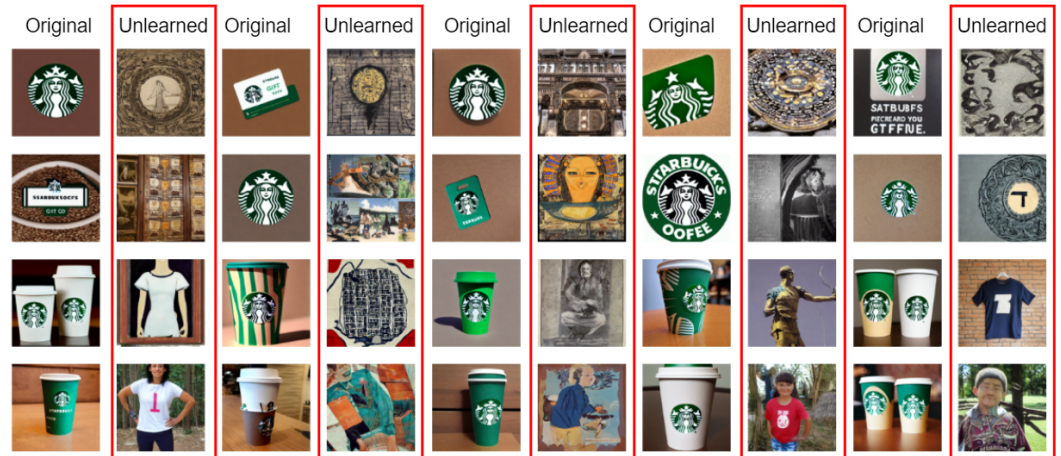


Figure 11: Unlearning results after erasing the concept of *Starbucks* logo (JS method).

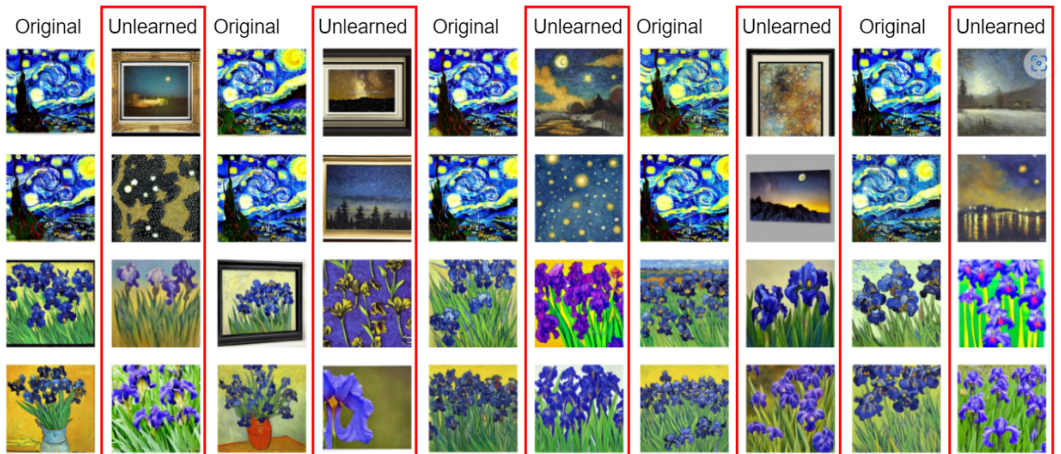


Figure 12: Unlearning results after erasing the concept of *Vincent Van Gogh* (JS method).

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

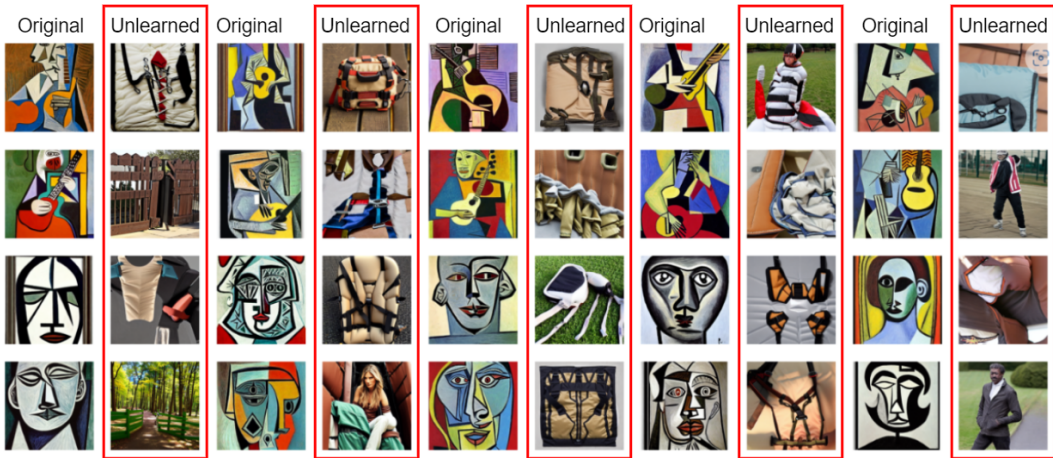


Figure 13: Unlearning results after erasing the concept of *Pablo Picasso* (JS method).



Figure 14: Unlearning results after erasing the concept of *Nike* logo (JS method).

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

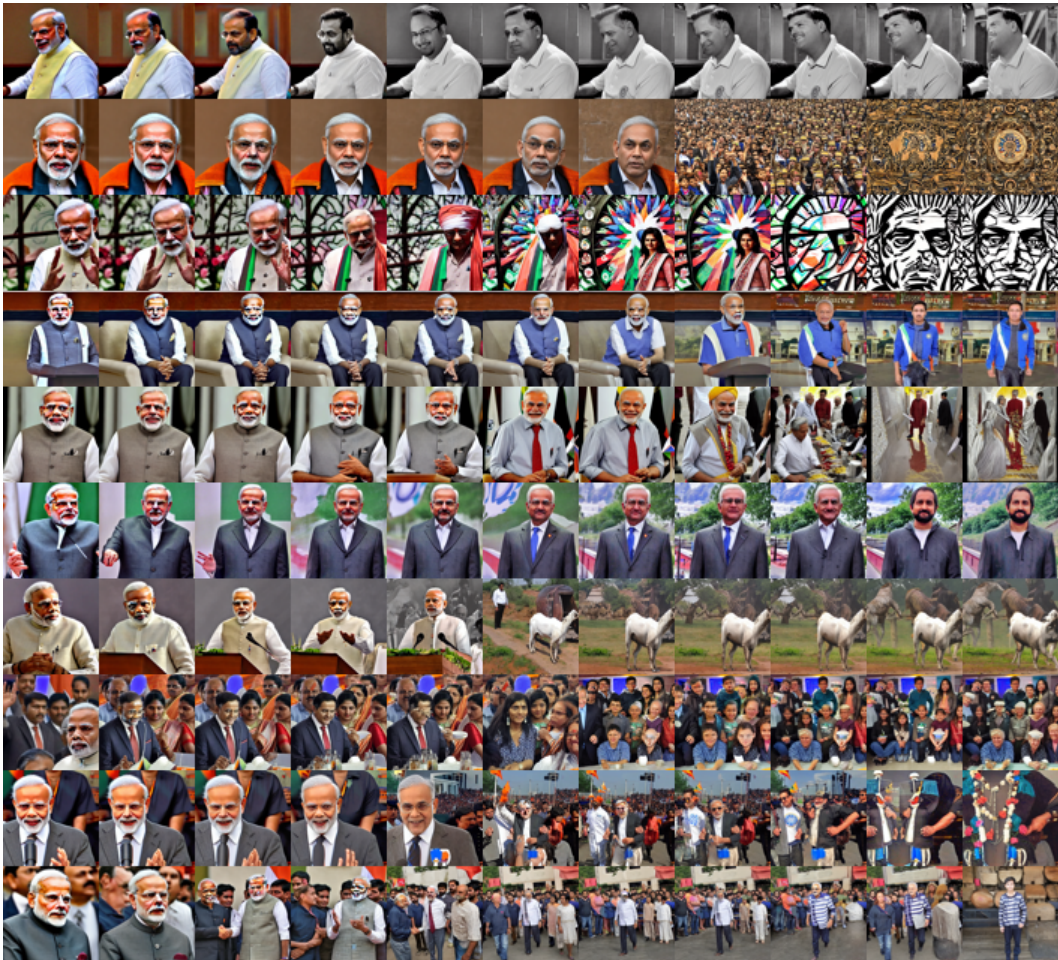


Figure 15: Unlearning results after erasing the concept of *Narendra Modi* (JS method).

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
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 16: Unlearning results after erasing the concept of *Amitabh Bachchan* (JS method).

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079



Figure 18: Unlearning results after erasing the concept of *Child* (JS method).

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

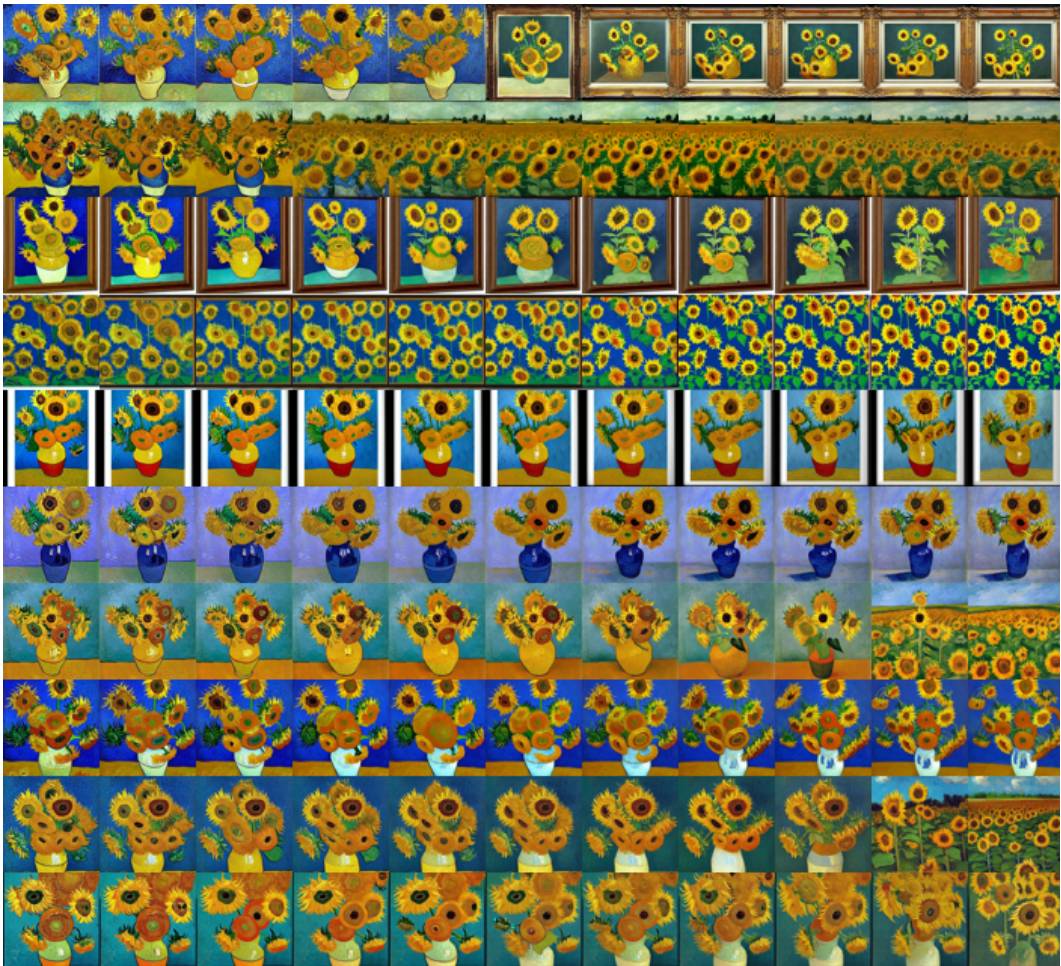


Figure 19: Unlearning results after erasing the concept of *Vincent Van Gogh* (JS method).

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187



Figure 20: Unlearning results after erasing the concept of *Pablo Picasso* (JS method).