

---

# UNLEARNCANVAS: A Stylized Image Dataset for Enhanced Machine Unlearning Evaluation in Diffusion Models

---

Yihua Zhang<sup>1</sup>, Chongyu Fan<sup>1</sup>, Yimeng Zhang<sup>1</sup>, Yuguang Yao<sup>1</sup>, Jinghan Jia<sup>1</sup>, Jiancheng Liu<sup>1</sup>, Gaoyuan Zhang<sup>2</sup>, Gaowen Liu<sup>3</sup>, Ramana Kompella<sup>3</sup>, Xiaoming Liu<sup>1</sup>, Sijia Liu<sup>1,2</sup>  
<sup>1</sup>Michigan State University, <sup>2</sup>IBM Research, <sup>3</sup>Cisco Research

## Abstract

The technological advancements in diffusion models (DMs) have demonstrated unprecedented capabilities in text-to-image generation and are widely used in diverse applications. However, they have also raised significant societal concerns, such as the generation of harmful content and copyright disputes. Machine unlearning (MU) has emerged as a promising solution, capable of removing undesired generative capabilities from DMs. However, existing MU evaluation systems present several key challenges that can result in incomplete and inaccurate assessments. To address these issues, we propose UNLEARNCANVAS, a comprehensive high-resolution stylized image dataset that facilitates the evaluation of the unlearning of artistic styles and associated objects. This dataset enables the establishment of a standardized, automated evaluation framework with 7 quantitative metrics assessing various aspects of the unlearning performance for DMs. Through extensive experiments, we benchmark 9 state-of-the-art MU methods for DMs, revealing novel insights into their strengths, weaknesses, and underlying mechanisms. Additionally, we explore challenging unlearning scenarios for DMs to evaluate worst-case performance against adversarial prompts, the unlearning of finer-scale concepts, and sequential unlearning. We hope that this study can pave the way for developing more effective, accurate, and robust DM unlearning methods, ensuring safer and more ethical applications of DMs in the future. The dataset, benchmark, and codes are publicly available at <https://unlearn-canvas.netlify.app/>.

## 1 Introduction

The recent technological breakthroughs in text-to-image generation, driven by diffusion models (DMs), have shown an unprecedented capability to produce high-resolution, high-quality images across a diverse range of subjects [1–10]. These models have become widely accessible to the public and are applied across various sectors, including advertising, the creative arts [11], and forensic sketching [12]. One reason for DMs being able to generate a broad spectrum of content is their reliance on diverse internet-sourced data [13]. However, this inclusiveness comes with risks, such as harmful generation [14], copyright issues [15], and biases or stereotypes [16, 17].

To alleviate the negative social impacts associated with DMs, machine unlearning (MU) techniques are catching increasing attention and have been studied in the field of text-to-image generation via DMs [18–33]. MU for DMs, also referred to as *DM unlearning* [27], aims to prevent the model from generating images when conditioned on an undesired concept (typically specified by a text prompt to be forgotten). Therefore, the problem of DM unlearning is also known as *concept erasing* [18–26]. Based on the types of unlearning requests, the targeted concept to be erased can be diverse, including not-safe-for-work (NSFW) prompts [14, 23], concrete object entities [18–33], and copyrighted information like artistic styles [23–26, 28–30, 32]. In Fig. 1 (a), we provide an illustration of DM

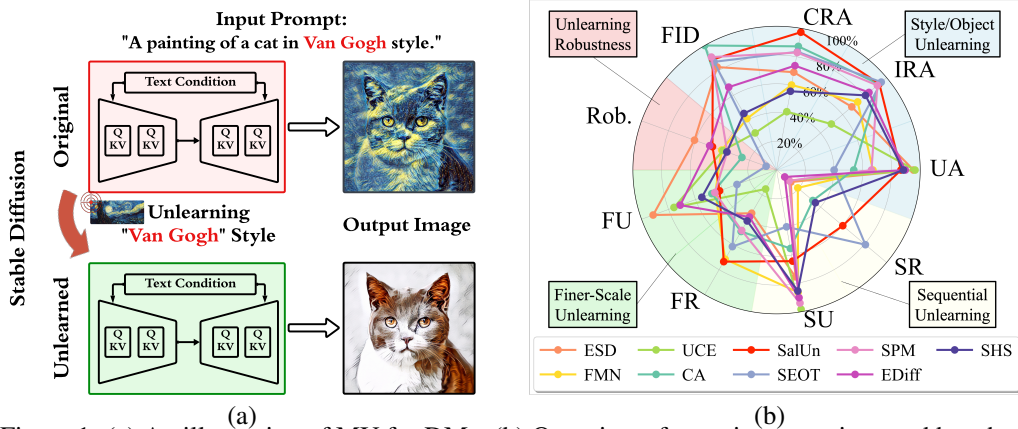


Figure 1: (a) An illustration of MU for DMs. (b) Overview of experiment settings and benchmark results. This benchmark focuses on three categories of quantitative metrics: the unlearning effectiveness (UA, Rob., FU, SU); the retainability of innocent knowledge (IRA, CRA, FR, SR); and the image generation quality (FID). Results are normalized to 0% ~ 100% per metric. No single method excels across all metrics. See a summary of these metrics in Tab. A1 and more results in Sec. 4.

unlearning by demonstrating the removal of an undesired artistic style from DM-generated images. Additionally, we refer readers to Sec. 2 for more related work on DM unlearning.

Compared to MU for DMs, the unlearning studies in computer vision have primarily focused on image classification models [34–48]. Given an image classification dataset, benchmarking unlearning performance for discriminative models is simpler, because MU typically involves either class-wise forgetting [49, 35] or data-wise forgetting [27, 40]. The former aims to eliminate the influence of a specific image class, while the latter removes the influence of specific data points from the entire training set. In contrast, DM unlearning is often applied to scenarios involving the removal of higher-level and more abstract concepts, rather than specific training data points, which can also be difficult to localize within the DM training set. The primary method for assessing the performance of DM unlearning is using the I2P (inappropriate image prompts) dataset [14]. This dataset mainly focuses on the safety assessment of *unlearned DMs* (i.e., DMs post-unlearning), designed to erase specific harmful concepts like ‘nudity’ and ‘violence’. Although I2P specifies which concepts to unlearn, evaluating the unlearning effectiveness and generation retainability of unlearned DMs under I2P remains challenging. This difficulty arises primarily from the lack of ground-truth data or objective criteria to precisely define safety versus non-safety. Furthermore, I2P cannot evaluate the performance of DM unlearning in other scenarios, such as object unlearning [19–33] and style unlearning [23–30].

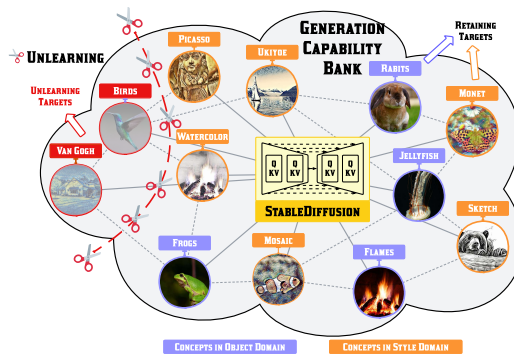


Figure 2: An illustration of machine unlearning using UNLEARNCANVAS. Concepts in the knowledge bank are categorized into different domains (style and object) and serve as potential unlearning targets. When one concept is unlearned, the rest concepts in both the same and different domains are required to be retained.

To enhance the assessment of machine unlearning in DMs and establish a standardized evaluation framework, we propose the development of a new benchmark dataset, referred to as UNLEARNCANVAS. Unlike I2P, UNLEARNCANVAS is designed to evaluate the unlearning of artistic painting styles along with associated image objects (see an illustration in Fig. 2). It also provides ground-truth data annotations to precisely define the criteria for the effectiveness of unlearning and the retained model utility post-unlearning (see a result overview in Fig. 1 (b)). We summarize **our contributions** below.

- We conduct a systematic review of existing MU methods for DMs and identified three unresolved challenges in their evaluations. We further provide an in-depth analysis of why a standardized benchmark for DM unlearning evaluation is crucial.

- We propose UNLEARNCANVAS, a large-scale, high-resolution stylized image dataset with diverse styles and objects, to address current challenges in DM unlearning evaluation. Its dual supervision of styles and objects ensures stylistic consistency, enabling a comprehensive, standardized evaluation approach to precisely characterize unlearning efficacy, generation utility preservation, and efficiency.
- We benchmark 9 state-of-the-art DM unlearning methods using UNLEARNCANVAS, covering standard assessments and examining more challenging scenarios, such as adversarial robustness, grouped object-style unlearning, and sequential unlearning. This comprehensive evaluation provides previously unknown insights into their strengths and weaknesses.

## 2 Related Work, Problem Statement, and Open Challenges

**Related work.** MU (machine unlearning) was originally developed to mitigate the (potentially detrimental) influence of specific data points in a pretrained ML model, without necessitating a complete retraining of the model after removing these unlearning data [34–52]. The significance of MU has emerged with the purpose of data privacy protection, in response to regulations such as ‘the right to be forgotten’ [53]. However, it has rapidly gained recognition for its role in promoting security, safety, and trustworthiness of ML models. Examples include defense against backdoor attacks [38, 54], fairness enhancement [55, 56], and controllable federated learning [57, 58]. More recently, DM unlearning has proven crucial in text-to-image generation to prevent the production of harmful, private, or illegal content [23–32]. For example, DMs have demonstrated a susceptibility to generating NSFW (not-safe-for-work) images when conditioned on inappropriate text prompts (*e.g.*, ‘nudity’ and ‘violence’) [14, 59, 60]. Similarly, MU has also found applications in preventing the generation of copyrighted information, such as the misuse of artistic painting styles [23, 28]. In brief, DM unlearning can be regarded as a model edit operation, aimed at preventing DMs from producing undesired or inappropriate images.

As unlearning approaches continue to evolve, the need for *benchmarking* their performance becomes increasingly critical. Standard datasets like CIFAR-10 [61] and ImageNet [62] have long been used to evaluate unlearning performance when forgetting a specific subset of data points [27, 40] or a targeted image class [49, 35]. Additionally, a recent benchmark specialized for unlearning in face recognition has been developed [63]. However, crafting effective benchmarks for generative models introduces additional challenges. In the language domain, recent initiatives have specifically aimed to evaluate the efficacy of MU for large language models (LLMs). Notable examples include TOFU [64], which explores the unlearning of fictitious elements, and WMDP [65], aimed at removing hazardous knowledge from LLMs. In contrast, to the best of our knowledge, there is currently no effective benchmark dataset designed for DM unlearning in text-to-image generation. The predominant dataset, I2P [14], targets the removal of harmful concepts for safe image generation. However, I2P only covers a specialized unlearning scenario and lacks precise evaluation criteria to support a comprehensive, quantitative, and precise assessment of unlearning performance. This gap inspires us to develop a new benchmark dataset for DM unlearning.

**Problem statement.** Let us denote a pretrained DM by  $\theta_o$ , capable of synthesizing high-quality images conditioned on a provided text prompt  $c$  (*e.g.*, ‘A painting of a cat in Van Gogh style’). When there is a request to adapt the model  $\theta_o$  to prevent generating images in a specific concept, such as the ‘Van Gogh’ painting style, the MU problem arises. Here, the ‘Van Gogh’ style becomes the specified *unlearning target*, also known as the *erasing concept* [23]. See Fig. 1-(a) for a visual representation of this unlearning problem. In this work, we specify the unlearning target as artistic styles and/or object classes to encompass the settings of style unlearning and object unlearning as described in the literature [23–32].

Table 1: Overview of existing MU evaluations for DMs, which covers two unlearning scenarios: object unlearning and style unlearning. Each case is demonstrated by the number of targeted unlearning concepts (target #), qualitative evaluation (Qual.) via generated image visualization, and quantitative evaluation (Quant.) of unlearning effectiveness (UE) and retainability (RT).

Related Work	Object Unlearning				Style Unlearning			
	Target #	Quant. UE	RT	Qual.	Target #	Quant. UE	RT	Qual.
ESD [23]	11	✓	✓	✓	5	✗	✗	✓
CA [16]	4	✓	✗	✓	4	✗	✗	✓
FMN [28]	7	✓	✗	✓	1	✗	✗	✓
UCE [24]	11	✓	✗	✓	7	✓	✗	✓
SA [29]	31	✓	✗	✓	1	✗	✗	✓
SalUn [27]	11	✓	✓	✓	0	✗	✗	✗
SPM [26]	24	✓	✓	✓	5	✗	✗	✓
SEOT [30]	22	✓	✓	✓	2	✗	✗	✓
EDiff [31]	2	✓	✓	✓	0	✗	✗	✗
SHS [32]	1	✓	✗	✓	0	✗	✗	✗

**Motivation: Unresolved challenges in MU evaluation.** When building a rigorous *quantitative* evaluation framework for DM unlearning, we identified three urgent challenges (C1)-(C3) after carefully examining existing studies as shown in **Tab. 1**.

(C1) *The absence of a consensus on a diverse unlearning target test repository.* As shown in **Tab. 1**, assessments of DM unlearning in various studies, in terms of both effectiveness (*i.e.*, concept erasure performance) and retainability (*i.e.*, preserved generation quality under non-forgotten concepts), typically use manually-selected unlearning targets from a *limited pool* and focus on *object-centric* evaluation. Precise evaluation for style unlearning and grouped object-style unlearning is lacking. Although the existing dataset WIKIART [66] features a collection of real-world artworks by various artists, it does *not* apply to object unlearning.

(C2) *The lack of a systematic study on ‘retainability’ of DMs post-unlearning.* As highlighted in the ‘retainability’ column of **Tab. 1**, there is a notable deficiency in the quantitative assessment of retainability for DMs post-unlearning. Retainability is crucial for measuring the potential side effects of a MU method. For instance, when unlearning an artistic style, it is essential to assess the DM’s ability to generate images in other styles and all objects. This makes the WIKIART dataset less suitable for benchmarking DM unlearning, as it lacks object labels needed to characterize the side effects of style unlearning on object recognition.

(C3) *The precision challenge in evaluating DM-generated images.* Artistic styles are inherently complex to define and differentiate precisely, complicating the quantitative evaluation of forgetting and retaining performance in style unlearning. For example, the WIKIART dataset was recently used in [67] to train a classifier for detecting copyright infringement, merely achieving an accuracy of 72.80%. The difficulty in precisely recognizing artistic styles in WIKIART may further hamper the evaluation of DMs’ post-unlearning generation. **Tab. A3** illustrates the challenge of accurately recognizing the styles of images generated by DMs, even when finetuned on WIKIART. This is evidenced by the test-time classification accuracy using the style classifier, ViT-L/16 [68], finetuned on WIKIART. As we can see, only about half of the images are correctly classified into their corresponding style classes, *e.g.*, 56.7% accuracy for classification on generated images using finetuned stable diffusion (SD) v2.0 [1], which is significantly lower than that on test images from the original WIKIART (85.4%).

The above challenges (C1)-(C3) drive us to develop the UNLEARNCANVAS dataset as a solution for the systematic and comprehensive evaluation of DM unlearning in **Sec. 3**.

### 3 Our Proposal: UNLEARNCANVAS Dataset

**Construction of UNLEARNCANVAS.** As motivated in **Sec. 2**, UNLEARNCANVAS is created for ease of MU evaluation in DMs. Its construction involves two main steps: seed image collection and subsequent image stylization; see **Fig. 3** for a schematic overview. More information about the dataset is provided in **Appx. A**.

*Seed image collection.* UNLEARNCANVAS includes images in 60 unique artistic styles across 20 distinct objects. The style categories are built upon a set of high-resolution *seed images*, given by real-world photos from Pexels [69]. There are 20 seed images collected for each of the 20 object classes. See step 1 in **Fig. 3**.

*Image stylization.* After collecting the seed images, the next step is the *stylization* process. During this phase, each seed image is transformed into 60 predetermined artistic styles, ensuring high stylistic consistency while preserving the original content details. This stylization process is carried out using services provided by Fotor [11]. Once all seed images have been stylized in all 60 styles, the dataset is constructed with super-high-resolution images and organized in a hierarchical structure that balances both style and object categories. See step 2 in **Fig. 3**.

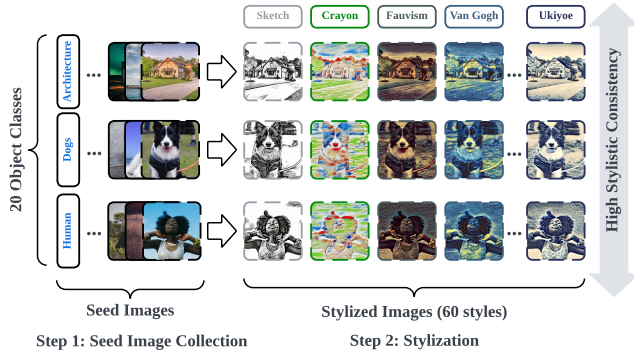


Figure 3: Illustration of curating UNLEARNCANVAS.

**Advantages of UNLEARNCANVAS.** Next, we will elaborate on the advantages (A1)-(A3) of UNLEARNCANVAS to address the MU evaluation challenges (C1)-(C3) discussed in Sec. 2.

(A1) *Style-object dual supervision enables a rich unlearning target bank for comprehensive evaluation.* To address (C1), UNLEARNCANVAS offers a diverse range of unlearning targets by annotating each image with both style and object labels. We categorized the unlearning targets into two main domains: a style domain comprising 60 styles, and an object domain with 20 object classes. This setup allows us to investigate the impact of style unlearning on the fidelity of object generation and vice versa.

(A2) *Enabling both ‘in-domain’ and ‘cross-domain’ retainability analyses for MU evaluation.* To address (C2), UNLEARNCANVAS facilitates us to assess DMs’ retainability post-unlearning from two perspectives: *in-domain retainability*, which evaluates the DM’s generation performance within the same domain as the unlearning target, and *cross-domain retainability*, which assesses the model’s generation ability under conditions from a different domain. For example, as depicted in Fig. 4, we can define the ‘domain’ as either ‘artistic style’ or ‘object class’. If the unlearning target is the ‘Van Gogh’ style, then generating images in all styles except ‘Van Gogh’ within the same object class (e.g., ‘Dog’ in Fig. 4) represents the in-domain case. Conversely, cross-domain retainability refers to generating images of different object classes (e.g., ‘Horses’ in Fig. 4), corresponding to an object domain different from the unlearning target within the style domain.

(A3) *High stylistic consistency ensures precise style definitions and enables accurate quantitative evaluations.* In UNLEARNCANVAS, images labeled with the same style exhibit high intra-style consistency and distinct inter-style differences. This facilitates the distinction between different styles, effectively addressing challenge (C2). In a similar study to that presented in Tab. A3, we observe that the styles within UNLEARNCANVAS can be readily classified using ViT-Large [68], achieving an accuracy of 98% for images generated by DM finetuned on UNLEARNCANVAS. For detailed results, please refer to Tab. A4.

The above advantages (A1)-(A3) make UNLEARNCANVAS uniquely suitable for assessing the performance of DM unlearning. For instance, UNLEARNCANVAS includes a greater number of high-resolution images (15M) compared to the existing dataset WIKIART (2M). Moreover, UNLEARNCANVAS surpasses WIKIART in terms of both intra-style coherence and inter-style distinctiveness. See Appx. C for more detailed comparisons.

**Evaluation pipeline via UNLEARNCANVAS.** We next introduce how UNLEARNCANVAS can be used to benchmark the performance of DM unlearning and the corresponding evaluation metrics. Fig. 5 illustrates the proposed evaluation pipeline with a designated unlearning target. It comprises four phases (I-IV) to evaluate unlearning effectiveness, retainability, generation quality, and efficiency.

*Phase I: Testbed preparation (Fig. 5 (a)).* Prior to unlearning, we commence by finetuning a specific DM, given by SD (StableDiffusion) v1.5, on UNLEARNCANVAS for text-to-image generation. We also train a vision transformer (ViT-Large) [68] for the purposes of style and object classification after unlearning. We ensure that both the SD model and classifiers attain high performance in image generation and classification across all styles and objects. See detailed testbed preparation in Appx. B.

*Phase II: Machine unlearning (Fig. 5 (b)).* Next, we utilize a given MU method to update the DM acquired in Phase I, aiming to eliminate a designated unlearning target (e.g., Van Gogh style). This results in the unlearned DM (denoted as  $\theta_u$ ), which avoids generating images in response to Van Gogh style-specific prompts.

*Phase III: Answer set generation (Fig. 5 (c)).* Further, we utilize the unlearned model  $\theta_u$  to generate a set of images conditioned on both the unlearning-related prompts and other benign prompts. To comprehensively evaluate unlearning performance, three types of answer sets are generated. (a) To assess unlearning effectiveness, images are generated in response to prompts that involve the unlearning target (e.g., any prompt like ‘An image in Van Gogh style’ for the case in Fig. 5). (b)



Figure 4: Illustration of in-domain and cross-domain retainability evaluation, with the Van Gogh style as the unlearning target. ✓ and ✗ indicate satisfactory and undesired results post unlearning.

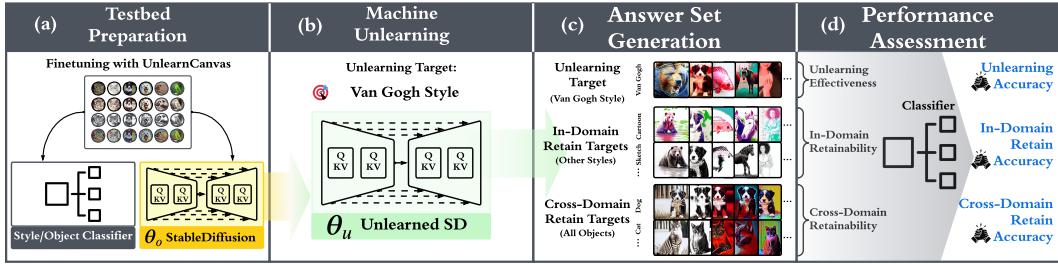


Figure 5: An illustration of the evaluation pipeline proposed in this work using UNLEARNCANVAS when unlearning a specific target concept ‘Van Gogh Style’. Unlearning performances (including the unlearning effectiveness and retainability) are quantitatively assessed (marked in blue) to accurately reflect the unlearning performance portrait. The unlearning target of the pipeline could traverse all the styles and objects to achieve a comprehensive evaluation.

To assess in-domain retainability (*i.e.*, the unlearned model’s ability to retain generation quality within the same domain), images are generated with prompts in the same domain as the unlearning target (*e.g.*, style-related prompts for style unlearning). These prompts cover all other styles provided in UNLEARNCANVAS except the one to be unlearned. (c) To measure cross-domain retainability, images are generated with prompts in other domains. For instance, in the case of style unlearning as shown in Fig. 5, the object domain is considered separate from the style domain. Prompts related to objects are used to evaluate the DM’s cross-domain generation capability after style unlearning.

*Phase IV: MU performance assessment (Fig. 5 (d)).* After Phase III, the answer set undergoes style/object classification for unlearning performance assessment. This classification results in three quantitative metrics ①-③. ① Unlearning accuracy (UA): This represents the proportion of images generated by the unlearned DM using the unlearning target-related prompt that is *not* correctly classified into the corresponding class. A *higher UA* indicates better unlearning performance in preventing image generation from the unlearning target-related prompts. ② In-domain retain accuracy (IRA): This is the classification accuracy of images generated by the unlearned DM using innocent prompts (not relevant to the unlearning target) within the same domain. ③ Cross-domain retain accuracy (CRA): Similar to IRA, this is the classification accuracy of images generated by the unlearned DM using innocent prompts in different domains. In addition to the accuracy metrics, we also use the **FID** score (④) to evaluate the distribution-wise generation quality of the unlearned DM. Furthermore, we monitor the efficiency of the unlearning process, considering factors such as **run-time** (⑤), **storage** space requirements (⑥), and **memory** costs (⑦).

**DM unlearning methods to be benchmarked.** In this work, we assess **9** most recently proposed MU methods for DMs, including **ESD** [23], **FMN** [28], **UCE** [24], **CA** [25], **SalUn** [27], **SEOT** [30], **SPM** [26], **EDiff** [31], and **SHS** [32]. We remark that SA [29], a recently proposed method, is excluded from our evaluation due to its excessive computational resource requirements and time consumption. Unless specified otherwise, SD v1.5 is the model used for unlearning. Detailed training settings for each method can be found in Appx. B.

## 4 Experiment Results

In this section, our benchmarking results are divided into two main parts. In Sec. 4.1, we comprehensively evaluate the existing DM unlearning methods on the established tasks of style and object unlearning. Extensive studies following the proposed evaluation pipeline (Fig. 5) reveal that focusing solely on unlearning effectiveness can lead to a biased perspective for DM unlearning if retainability is not concurrently assessed (Tab. 2 & Fig. A2). We also show that preserving CRA (cross-domain retainability) is more challenging than IRA (in-domain retainability), highlighting a gap in the current literature (Fig. 6). Furthermore, we find that different DM unlearning methods exhibit distinct unlearning mechanisms by examining their unlearning directions (Fig. 7). In Sec. 4.2, we introduce more challenging unlearning scenarios. First, we assess the robustness of current DM unlearning methods against adversarial prompts crafted based on (Fig. 8). Additionally, we examine the performance when facing unlearning targets with finer granularity, formed by style-object concept combinations (Tab. 3). Further, we leverage UNLEARNCANVAS to provide insights into DM unlearning in a sequential unlearning fashion (Tab. A5).

Table 2: Performance overview of different DM unlearning methods evaluated on UNLEARNCANVAS. The performance metrics include UA (unlearning accuracy), IRA, CRA, and FID. The symbols  $\uparrow$  or  $\downarrow$  indicate whether a higher or lower value represents better performance. Results are averaged over all the style and object unlearning cases. The best performance for each metric is highlighted in green, while significantly underperforming results are marked in red, indicating areas needing improvement for existing DM unlearning methods.

Method	Effectiveness						Efficiency			
	Style Unlearning			Object Unlearning			FID ( $\downarrow$ )	Time (s) ( $\downarrow$ )	Memory (GB) ( $\downarrow$ )	Storage (GB) ( $\downarrow$ )
UA ( $\uparrow$ )	IRA ( $\uparrow$ )	CRA ( $\uparrow$ )	UA ( $\uparrow$ )	IRA ( $\uparrow$ )	CRA ( $\uparrow$ )					
ESD [23]	98.58%	80.97%	93.96%	92.15%	55.78%	44.23%	65.55	6163	17.8	4.3
FMN [28]	88.48%	56.77%	46.60%	45.64%	90.63%	73.46%	131.37	350	17.9	4.2
UCE [24]	98.40%	60.22%	47.71%	94.31%	39.35%	34.67%	182.01	434	5.1	1.7
CA [25]	60.82%	96.01%	92.70%	46.67%	90.11%	81.97%	54.21	734	10.1	4.2
SalUn [27]	86.26%	90.39%	95.08%	86.91%	96.35%	99.59%	61.05	667	30.8	4.0
SEOT [30]	56.90%	94.68%	84.31%	23.25%	95.57%	82.71%	62.38	95	7.34	0.0
SPM [26]	60.94%	92.39%	84.33%	71.25%	90.79%	81.65%	59.79	29700	6.9	0.0
EDiff [31]	92.42%	73.91%	98.93%	86.67%	94.03%	48.48%	81.42	1567	27.8	4.0
SHS [32]	95.84%	80.42%	43.27%	80.73%	81.15%	67.99%	119.34	1223	31.2	4.0

#### 4.1 Benchmarking Current DM Unlearning Methods for Style and Object Unlearning

**Overall performance of DM unlearning.** In Tab. 2, we provide an overview of the performance of existing unlearning methods, using the evaluation metrics associated with UNLEARNCANVAS.

First, as seen in Tab. 2, retainability is essential for a comprehensive assessment of DM unlearning. For example, in the scenario of style unlearning, ESD and UCE achieve similar UA (both around 98%), but their retainability (IRA and CRA) differs significantly, with one over 80% and the other around 60%. Therefore, relying solely on UA can provide a skewed view of the performance of DM unlearning. Second, we observe that CRA is harder to retain than IRA. For example, UCE exhibits a significant gap between its IRA and CRA, both in style (60.22% vs. 47.71%) and object (39.35% vs. 34.67%) unlearning scenarios. Similar patterns are observed with other methods like FMN, CA, SEOT, SPM, and SHS. This pattern suggests that while MU methods are somewhat effective at preserving concepts within the same domain, they face greater challenges in maintaining performance on unlearning target-unrelated concepts across domains. This issue has been overlooked in previous studies due to the absence of a multi-label dataset like UNLEARNCANVAS for systematic unlearning analyses. Third, we find that a single unlearning method can perform differently across various domains, and no single method excels in all aspects. For example, FMN achieves a UA of 88.48% in style unlearning but a much lower UA of 45.64% in object unlearning. A similar phenomenon can be observed with the unlearning methods SEOT, SPM, and SHS, which exhibit significantly large performance gaps between style and object unlearning. Furthermore, each method exhibits unique strengths and weaknesses. For example, while ESD typically shows improvement in UA, it may lag in terms of IRA and CRA, particularly in the case of object unlearning. Some methods, such as SEOT and SPM, achieve extremely high storage efficiency, requiring no additional storage, but suffer from inferior UA (below 80%). These observations underscore the challenges of DM unlearning and highlight the need for further advancements in the field.

**A closer look into retainability: A case study on ESD.** We next perform an in-depth analysis of DMs’ retainability after unlearning, focusing on the most popular ESD method [23]. In Fig. 6 (left), we present a heatmap illustrating the per-style/object unlearning accuracy (diagonal) and retain accuracy (off-diagonal) for ESD. The heatmap segments are labeled to indicate different regions of interest (ROIs) in style and object unlearning related to our evaluation metrics: **A1/A2** for UA, **B1/B2** for IRA, and **C1/C2** for CRA. Key observations are highlighted below.

First, we observe distinct behaviors of ESD in unlearning styles compared to objects. For example, ROI C2 appears darker than C1, suggesting a reduced ability to retain styles during object unlearning. Similar patterns are seen between B2 and B1. For a clearer explanation, Fig. 6 (right) provides image generation examples of the DM before and after unlearning within different ROIs. Second, We observe that style/object unlearning is relatively easier compared to retaining the generation performance of unlearned DMs conditioned on unlearning-unrelated prompts. This is evident in the contrast between ROIs A1/A2 and ROIs B1/B2 or C1/C2. One exception is the case of unlearning the object ‘Sea’, which exhibits relatively low UA in A2. This may be attributed to post-unlearning image generation still resembling sea waves, as shown by the image examples in Fig. 6 (right) for

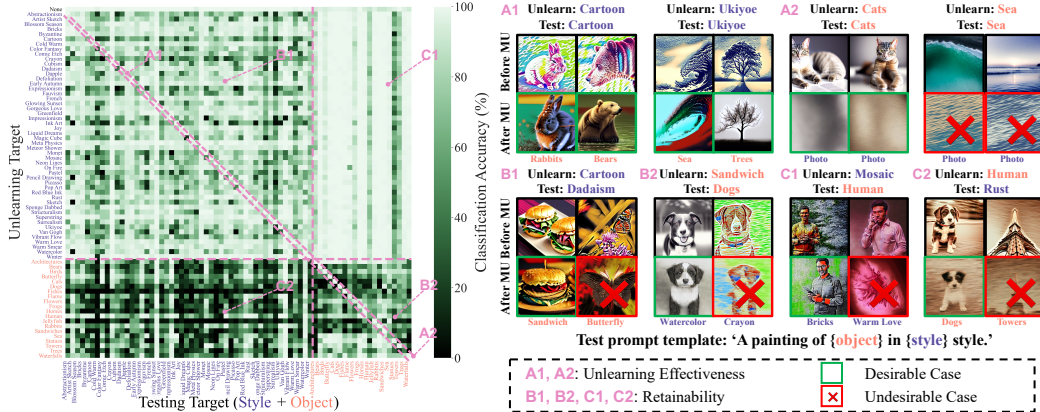


Figure 6: *Left*: Heatmap visualization of the unlearning accuracy and retainability of ESD on UNLEARNCANVAS. The  $x$ -axis shows the tested concepts for image generation using the unlearned model, while the  $y$ -axis indicates the unlearning target. Concept types are distinguished by color: **styles** in blue and **objects** in orange. The figure is divided into regions representing corresponding evaluation metrics and unlearning scopes ( $A$  for UA,  $B$  for IRA,  $C$  for CRA; ‘1’ for style unlearning, ‘2’ for object unlearning). Higher values in lighter colors denote better performance. The first row serves as a reference for comparison before unlearning. Zooming into the figure is recommended for detailed observation. *Right*: Representative cases illustrating each region with images generated before and after unlearning a specific concept.

A2. Third, we observe an inherent trade-off between unlearning effectiveness and retainability. For example, while SalUn exhibits more consistent performance across different scenarios than ESD, it does not achieve the same level of UA, as indicated by the lighter diagonal values in Fig. A3. For the heatmap performance of other unlearning methods, please refer to Appx. D.

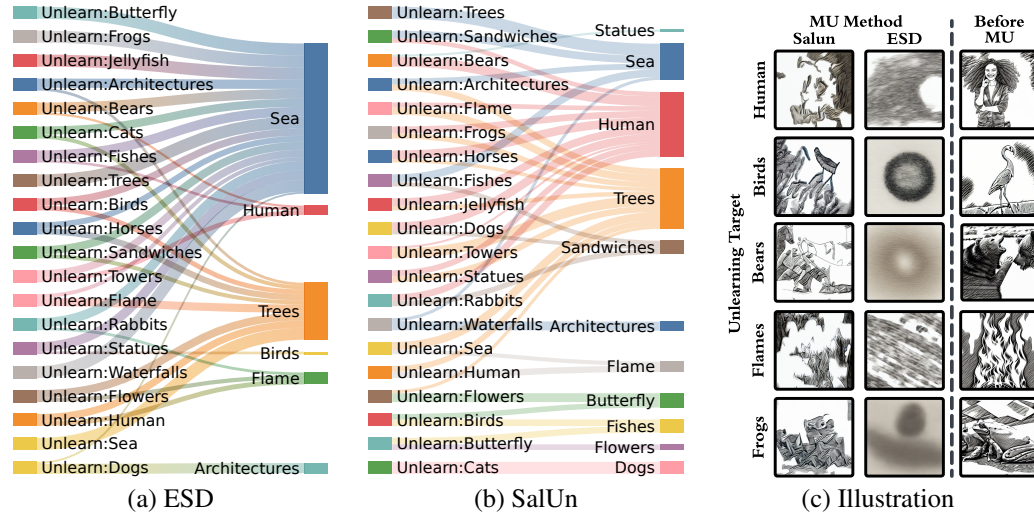


Figure 7: Visualization of the unlearning directions of (a) ESD and (b) SalUn. This figure illustrates the conceptual shift of the generated images of an unlearned model conditioned on the unlearning target. Images generated by the post-unlearning models are classified and used to understand this shift. Edges leading from the object in the left column to the right signify that images generated conditioned on unlearning targets are instead classified as the shifted concepts after unlearning. This reveals the primary unlearning direction for each unlearning method. The most dominant unlearning direction for an object is visualized. Figure (c) provides visualizations of generated images using the prompt template ‘A painting of {object} in Sketch style.’ with *object* being each unlearning target.

**Understanding unlearning method’s behavior via unlearning directions.** As noted earlier, different unlearning methods display distinct unlearning behaviors. To gain insights into the underlying reasons for these differences, Fig. 7 (a) and (b) visualize the ‘unlearning directions’ for ESD and



SalUn, respectively. These unlearning directions are determined by connecting the unlearning target with the predicted label of the generated image from the unlearned DM conditioned on the unlearning target. As shown in Fig. 7 (a), ESD demonstrates a focused shift in image generation after object unlearning, with a predominant transition towards generating images labeled by ‘Sea’ and ‘Trees’. This behavior arises from ESD’s optimization process, designed to steer the generation of the DM away from a predefined concept. Consequently, images generated by the ESD-induced unlearned model consistently lack clearly identifiable objects, resembling waves and trees, which leads to their classification into the ‘Sea’ and ‘Trees’ classes; see Fig. 7 (c) for examples of generated images. In contrast, SalUn exhibits a more diverse range of unlearning directions, shifting images to 11 different objects. This diversity results from SalUn’s requirement to replace the unlearning target with a random concept. As shown in Fig. 7 (c), images generated by SalUn post-object unlearning still maintain some object contours (different from the original unlearning target) and better retain style information compared to ESD.

## 4.2 Benchmarking Current DM Unlearning Methods in More Challenging Scenarios

### Unlearning robustness against adversarial prompts.

Recent studies [59, 60] have highlighted the vulnerabilities of DM unlearning to adversarial prompts, such as jailbreak attacks. In Fig. 8, we use the state-of-the-art attacking method, UnlearnDiffAtk [59], to craft adversarial prompts based on UNLEARNCANVAS and evaluate the robustness of unlearned DMs shown in Tab. 2. See Appx. B.4 for detailed setup. As we can see, all the DM unlearning methods experience a significant drop in UA, falling below 60% when confronted with adversarial prompts. Notably, methods like UCE, SEOT, and SHS exhibit particularly steep declines. Moreover, a high UA under normal conditions does not necessarily imply robustness to adversarial attacks. For instance, although UCE achieves a UA over 90% in normal settings, its performance plummets to below 50% against adversarial prompts. This stark contrast highlights the importance of the worst-case evaluation for MU methods.

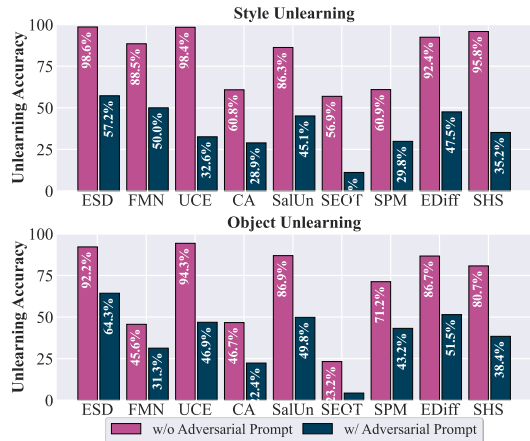


Figure 8: UA of DM unlearning against adversarial prompts [59]. Unlearned models in Tab. 2 are used as victim models to generate adversarial prompts.

### Unlearning at a finer scale: Style-object combinations as unlearning targets.

Next, we consider the fine granularity of the unlearning target, defined by a style-object combination, such as ‘‘An image of dogs in Van Gogh style’’. This unlearning challenge requires the unlearned model to avoid affecting image generation for dogs in non-Van Gogh styles and Van Gogh-style images with non-dog objects. See Appx. B.5 for detailed setup. Here, besides UA, the retainability performance is assessed in three contexts. (a) Style consistency (SC): retainability under prompts with different objects but in the same style, e.g., ‘‘An image of cats in Van Gogh style’’. (b) Object consistency (OC): retainability under prompts featuring the same object in different styles, e.g., ‘‘An image of dogs in Picasso style’’. (c) Unrelated prompting (UP): retainability for all other prompts. Tab. 2 presents the performance of unlearning style-object combinations. Style-object combinations are more challenging to unlearn than individual objects or styles, as evidenced by a significant drop in UA—over 20% lower compared to values in Tab. 2. Retainability drops to below 20% for top-performing methods like ESD and UCE, originally highlighted for their efficacy. This is presumably

Table 3: Performance of unlearning style-object combinations. The assessment includes UA and retainability in three contexts: SC (style consistency), OC (object consistency), and UP (unrelated prompting). The best result in each metric is highlighted in green.

Method	UA	SC	OC	UP
ESD [23]	91.42%	4.88%	14.72%	84.38%
FMN [28]	45.37%	68.73%	62.74%	83.25%
UCE [24]	75.97%	4.53%	5.72%	35.42%
CA [25]	47.92%	10.08%	56.35%	81.54%
SalUn [27]	42.21%	62.45%	70.93%	87.28%
SEOT [30]	29.32%	45.31%	53.64%	85.45%
SPM [26]	45.72%	41.34%	36.32%	67.82%
EDiff [31]	71.33%	35.23%	26.32%	51.52%
SHS [32]	55.32%	14.34%	24.32%	83.95%

due to ESD’s underlying unlearning mechanism, which requires only a single prompt, resulting in a poor ability to precisely define the unlearning scope.

**Evaluation in sequential unlearning.** Furthermore, we evaluate the unlearning performance in sequential unlearning (SU) [70, 71], where unlearning requests arrive sequentially. This parallels the continual learning (CL) task, which requires models to not only unlearn new targets effectively but also maintain the unlearning of previous targets while retaining all other knowledge. Here, we consider unlearning 6 styles sequentially and the results are presented in **Tab. A5**. We remark that the method SEOT does not support sequential unlearning in its original implementation and thus is not included in **Tab. A5**. Our findings reveal significant insights. (1) *Degraded retainability*: Sequential unlearning requests generally degrade retainability across all methods, with RA values frequently dropping below the average levels previously seen in **Tab. 2**. Here RA is given by the average of IRA and CRA. (2) *Unlearning rebound effect*: Knowledge previously unlearned can be inadvertently reactivated by new unlearning requests. This is evidenced by decreasing UA values for earlier objectives as more unlearning tasks are introduced. This suggests that residual knowledge remains within the model and can be reactivated, aligning with findings from **Fig. 8**. This indicates the unlearned models by some MU methods do not essentially lose the generation ability of the unlearning target. (3) *Catastrophic retaining failure*: RA significantly drops at a certain request, exemplified by a sudden decrease in RA of UCE from 81.42% to 29.38% after the second request,  $\mathcal{T}_2$ . This indicates that the seemingly acceptable side effects generated by some unlearning methods will drastically modify the knowledge representations when accumulated. This experiment illuminates the complex dynamics of knowledge removal and retention within DMs and highlights the potential pitfalls of existing unlearning methods when faced with sequential unlearning tasks. The observation of the ‘unlearning rebound effect’ and ‘catastrophic retaining failure’ particularly emphasizes the need for a more nuanced understanding of how knowledge is managed within DMs.

**Visualizations.** We provide plenty of visualizations, illustrations, and qualitative results. These visualizations are intended to deepen the understanding of the effects of different MU methods and clearly illustrate the challenges identified in previous sections. Specifically, in **Fig. A11**, we provide abundant generation examples of all the 9 methods benchmarked in this work in a case study of unlearning the ‘Cartoon’ style. Both the successful and failure cases are demonstrated in the context of unlearning effectiveness, in-domain retainability, and cross-domain retainability. Besides, in **Fig. A12**, we provide visualizations for the effect of adversarial prompts enabled by UnlearnDiffAtk [59].

## 5 Conclusion, Discussion, and Limitation

In this paper, we systematically reviewed existing MU (machine unlearning) methods on DMs (diffusion models) and identified key challenges in their evaluation systems that could lead to incomplete, inaccurate, and biased assessments. In response, we propose UNLEARNCANVAS, a high-resolution stylized image dataset designed to facilitate comprehensive evaluation of MU methods. We also introduce novel systematic evaluation metrics. By benchmarking nine state-of-the-art MU methods, we reveal novel insights into their strengths and weaknesses and also deepen the understanding of their underlying mechanisms. Our analysis of three challenging unlearning tasks highlights significant shortcomings of current methods, including a lack of robustness, difficulties in fine-scale unlearning, and issues arising from sequential unlearning tasks. These findings also point to meaningful future research directions aimed at developing satisfactory and practical MU methods for real-world applications. Furthermore, we recognize limitations in our benchmark, such as its focus on specific Stable Diffusion models and the text-to-image task in DMs. We hope this work serves as a foundation for future research to broaden MU evaluations across a wider range of model architectures and tasks in DMs.

## Acknowledgement

The work of Y. Zhang, C. Fan, Y. Zhang, Y. Yao, J. Jia, J. Liu, and S. Liu was supported by the Cisco Research Award, the ARO Award W911NF2310343, the National Science Foundation (NSF) CISE Core Program Award IIS-2207052, and the NSF CAREER Award IIS-2338068, and the Amazon Research Award for AI in Information Security.

## Bibliography

- [1] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.
- [2] “Midjourney,” <https://www.midjourney.com/>, 2023, accessed: 2024-01-19.
- [3] M. Kang, J.-Y. Zhu, R. Zhang, J. Park, E. Shechtman, S. Paris, and T. Park, “Scaling up gans for text-to-image synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 10 124–10 134.
- [4] M. Tao, B.-K. Bao, H. Tang, and C. Xu, “Galip: Generative adversarial clips for text-to-image synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14 214–14 223.
- [5] M. Tao, H. Tang, F. Wu, X.-Y. Jing, B.-K. Bao, and C. Xu, “Df-gan: A simple and effective baseline for text-to-image synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 515–16 525.
- [6] B. Li, X. Qi, T. Lukasiewicz, and P. Torr, “Controllable text-to-image generation,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [7] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee, “Generative adversarial text to image synthesis,” in *International conference on machine learning*. PMLR, 2016, pp. 1060–1069.
- [8] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, “Learning to prompt for vision-language models,” *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2337–2348, 2022.
- [9] B. Kawar, S. Zada, O. Lang, O. Tov, H. Chang, T. Dekel, I. Mosseri, and M. Irani, “Imagic: Text-based real image editing with diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 6007–6017.
- [10] M. Zhang, Z. Cai, L. Pan, F. Hong, X. Guo, L. Yang, and Z. Liu, “Motiondiffuse: Text-driven human motion generation with diffusion model,” *arXiv preprint arXiv:2208.15001*, 2022.
- [11] “Fotor: Online photo editor for everyone.” <https://www.fotor.com/>, 2023.
- [12] “Forensic sketch artist at lablab’s openai event,” 2023.
- [13] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman *et al.*, “Laion-5b: An open large-scale dataset for training next generation image-text models,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 25 278–25 294, 2022.
- [14] P. Schramowski, M. Brack, B. Deiseroth, and K. Kersting, “Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22 522–22 531.
- [15] I. Krietzberg, “Users of midjourney text-to-image site claim issues with new update,” *The Street*, 2023.
- [16] Y. Vinker, A. Voynov, D. Cohen-Or, and A. Shamir, “Concept decomposition for visual exploration and inspiration,” *ACM Transactions on Graphics (TOG)*, vol. 42, no. 6, pp. 1–13, 2023.
- [17] V. Turk, “How ai reduces the world to stereotypes,” *Rest of World*, 2023. [Online]. Available: [https://restofworld.org/2023/ai-image-stereotypes/?utm\\_source=pocket-newtab-en-us](https://restofworld.org/2023/ai-image-stereotypes/?utm_source=pocket-newtab-en-us)
- [18] Y.-L. Tsai, C.-Y. Hsu, C. Xie, C.-H. Lin, J.-Y. Chen, B. Li, P.-Y. Chen, C.-M. Yu, and C.-Y. Huang, “Ring-a-bell! how reliable are concept removal methods for diffusion models?” *arXiv preprint arXiv:2310.10012*, 2023.

- [19] P. Dong, S. Guo, J. Wang, B. Wang, J. Zhang, and Z. Liu, “Towards test-time refusals via concept negation,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [20] S. Hong, J. Lee, and S. S. Woo, “All but one: Surgical concept erasing with model preservation in text-to-image diffusion models,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 19, 2024, pp. 21 143–21 151.
- [21] M. Zhao, L. Zhang, T. Zheng, Y. Kong, and B. Yin, “Separable multi-concept erasure from diffusion models,” *arXiv preprint arXiv:2402.05947*, 2024.
- [22] C.-P. Huang, K.-P. Chang, C.-T. Tsai, Y.-H. Lai, and Y.-C. F. Wang, “Receler: Reliable concept erasing of text-to-image diffusion models via lightweight erasers,” *arXiv preprint arXiv:2311.17717*, 2023.
- [23] R. Gandikota, J. Materzynska, J. Fiotto-Kaufman, and D. Bau, “Erasing concepts from diffusion models,” *arXiv preprint arXiv:2303.07345*, 2023.
- [24] R. Gandikota, H. Orgad, Y. Belinkov, J. Materzyńska, and D. Bau, “Unified concept editing in diffusion models,” *arXiv preprint arXiv:2308.14761*, 2023.
- [25] N. Kumari, B. Zhang, S.-Y. Wang, E. Shechtman, R. Zhang, and J.-Y. Zhu, “Ablating concepts in text-to-image diffusion models,” *arXiv preprint arXiv:2303.13516*, 2023.
- [26] M. Lyu, Y. Yang, H. Hong, H. Chen, X. Jin, Y. He, H. Xue, J. Han, and G. Ding, “One-dimensional adapter to rule them all: Concepts, diffusion models and erasing applications,” *arXiv preprint arXiv:2312.16145*, 2023.
- [27] C. Fan, J. Liu, Y. Zhang, D. Wei, E. Wong, and S. Liu, “Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation,” *arXiv preprint arXiv:2310.12508*, 2023.
- [28] E. Zhang, K. Wang, X. Xu, Z. Wang, and H. Shi, “Forget-me-not: Learning to forget in text-to-image diffusion models,” *arXiv preprint arXiv:2303.17591*, 2023.
- [29] A. Heng and H. Soh, “Selective amnesia: A continual learning approach to forgetting in deep generative models,” *arXiv preprint arXiv:2305.10120*, 2023.
- [30] S. Li, J. van de Weijer, T. Hu, F. S. Khan, Q. Hou, Y. Wang, and J. Yang, “Get what you want, not what you don’t: Image content suppression for text-to-image diffusion models,” *arXiv preprint arXiv:2402.05375*, 2024.
- [31] J. Wu, T. Le, M. Hayat, and M. Harandi, “Erasediff: Erasing data influence in diffusion models,” *arXiv preprint arXiv:2401.05779*, 2024.
- [32] J. Wu and M. Harandi, “Scissorhands: Scrub data influence via connection sensitivity in networks,” *arXiv preprint arXiv:2401.06187*, 2024.
- [33] Y. Zhang, X. Chen, J. Jia, Y. Zhang, C. Fan, J. Liu, M. Hong, K. Ding, and S. Liu, “Defensive unlearning with adversarial training for robust concept erasure in diffusion models,” *arXiv preprint arXiv:2405.15234*, 2024.
- [34] Y. Cao and J. Yang, “Towards making systems forget with machine unlearning,” in *2015 IEEE symposium on security and privacy*. IEEE, 2015, pp. 463–480.
- [35] A. Goh, A. Achille, and S. Soatto, “Eternal sunshine of the spotless net: Selective forgetting in deep networks,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9304–9312.
- [36] A. Becker and T. Liebig, “Evaluating machine unlearning via epistemic uncertainty,” *arXiv preprint arXiv:2208.10836*, 2022.
- [37] A. Thudi, G. Deza, V. Chandrasekaran, and N. Papernot, “Unrolling sgd: Understanding factors influencing machine unlearning,” in *2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2022, pp. 303–319.

- [38] J. Jia, J. Liu, P. Ram, Y. Yao, G. Liu, Y. Liu, P. Sharma, and S. Liu, “Model sparsity can simplify machine unlearning,” in *Annual Conference on Neural Information Processing Systems*, 2023.
- [39] M. Chen, W. Gao, G. Liu, K. Peng, and C. Wang, “Boundary unlearning: Rapid forgetting of deep networks via shifting the decision boundary,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7766–7775.
- [40] A. Warnecke, L. Pirch, C. Wressnegger, and K. Rieck, “Machine unlearning of features and labels,” *arXiv preprint arXiv:2108.11577*, 2021.
- [41] C. Fan, J. Liu, A. Hero, and S. Liu, “Challenging forgets: Unveiling the worst-case forget sets in machine unlearning,” *arXiv preprint arXiv:2403.07362*, 2024.
- [42] A. Ginart, M. Guan, G. Valiant, and J. Y. Zou, “Making ai forget you: Data deletion in machine learning,” *Advances in neural information processing systems*, vol. 32, 2019.
- [43] S. Neel, A. Roth, and S. Sharifi-Malvajerdi, “Descent-to-delete: Gradient-based methods for machine unlearning,” in *Algorithmic Learning Theory*. PMLR, 2021, pp. 931–962.
- [44] E. Ullah, T. Mai, A. Rao, R. A. Rossi, and R. Arora, “Machine unlearning via algorithmic stability,” in *Conference on Learning Theory*. PMLR, 2021, pp. 4126–4142.
- [45] A. Sekhari, J. Acharya, G. Kamath, and A. T. Suresh, “Remember what you want to forget: Algorithms for machine unlearning,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 18 075–18 086, 2021.
- [46] L. Bourtole, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, and N. Papernot, “Machine unlearning,” in *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2021, pp. 141–159.
- [47] T. T. Nguyen, T. T. Huynh, P. L. Nguyen, A. W.-C. Liew, H. Yin, and Q. V. H. Nguyen, “A survey of machine unlearning,” *arXiv preprint arXiv:2209.02299*, 2022.
- [48] Z. Izzo, M. A. Smart, K. Chaudhuri, and J. Zou, “Approximate data deletion from machine learning models,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 2008–2016.
- [49] L. Graves, V. Nagisetty, and V. Ganesh, “Amnesiac machine learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 13, 2021, pp. 11 516–11 524.
- [50] C. Guo, T. Goldstein, A. Hannun, and L. Van Der Maaten, “Certified data removal from machine learning models,” *arXiv preprint arXiv:1911.03030*, 2019.
- [51] H. Xu, T. Zhu, L. Zhang, W. Zhou, and P. S. Yu, “Machine unlearning: A survey,” *ACM Computing Surveys*, vol. 56, no. 1, pp. 1–36, 2023.
- [52] J. Jia, Y. Zhang, Y. Zhang, J. Liu, B. Runwal, J. Diffenderfer, B. Kailkhura, and S. Liu, “Soul: Unlocking the power of second-order optimization for llm unlearning,” *arXiv preprint arXiv:2404.18239*, 2024.
- [53] C. J. Hoofnagle, B. van der Sloot, and F. Z. Borgesius, “The european union general data protection regulation: what it is and what it means,” *Information & Communications Technology Law*, vol. 28, no. 1, pp. 65–98, 2019.
- [54] Y. Liu, M. Fan, C. Chen, X. Liu, Z. Ma, L. Wang, and J. Ma, “Backdoor defense with machine unlearning,” *arXiv preprint arXiv:2201.09538*, 2022.
- [55] R. Chen, J. Yang, H. Xiong, J. Bai, T. Hu, J. Hao, Y. Feng, J. T. Zhou, J. Wu, and Z. Liu, “Fast model debias with machine unlearning,” *arXiv preprint arXiv:2310.12560*, 2023.
- [56] A. Oesterling, J. Ma, F. P. Calmon, and H. Lakkaraju, “Fair machine unlearning: Data removal while mitigating disparities,” *arXiv preprint arXiv:2307.14754*, 2023.
- [57] A. Halimi, S. Kadhe, A. Rawat, and N. Baracaldo, “Federated unlearning: How to efficiently erase a client in fl?” *arXiv preprint arXiv:2207.05521*, 2022.

- [58] R. Jin, M. Chen, Q. Zhang, and X. Li, “Forgettable federated linear learning with certified data removal,” *arXiv preprint arXiv:2306.02216*, 2023.
- [59] Y. Zhang, J. Jia, X. Chen, A. Chen, Y. Zhang, J. Liu, K. Ding, and S. Liu, “To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now,” *arXiv preprint arXiv:2310.11868*, 2023.
- [60] Z.-Y. Chin, C.-M. Jiang, C.-C. Huang, P.-Y. Chen, and W.-C. Chiu, “Prompting4debugging: Red-teaming text-to-image diffusion models by finding problematic prompts,” *arXiv preprint arXiv:2309.06135*, 2023.
- [61] A. Krizhevsky and G. Hinton, “Learning multiple layers of features from tiny images,” *Master’s thesis, Department of Computer Science, University of Toronto*, 2009.
- [62] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, 2009, pp. 248–255.
- [63] D. Choi and D. Na, “Towards machine unlearning benchmarks: Forgetting the personal identities in facial recognition systems,” *arXiv preprint arXiv:2311.02240*, 2023.
- [64] P. Maini, Z. Feng, A. Schwarzschild, Z. C. Lipton, and J. Z. Kolter, “Tofu: A task of fictitious unlearning for llms,” *arXiv preprint arXiv:2401.06121*, 2024.
- [65] N. Li, A. Pan, A. Gopal, S. Yue, D. Berrios, A. Gatti, J. D. Li, A.-K. Dombrowski, S. Goel, L. Phan *et al.*, “The wmdp benchmark: Measuring and reducing malicious use with unlearning,” *arXiv preprint arXiv:2403.03218*, 2024.
- [66] F. Phillips and B. Mackintosh, “Wiki art gallery, inc.: A case for critical thinking,” *Issues in Accounting Education*, vol. 26, no. 3, pp. 593–608, 2011.
- [67] M. Moayeri, S. Basu, S. Balasubramanian, P. Kattakinda, A. Chengini, R. Brauneis, and S. Feizi, “Rethinking artistic copyright infringements in the era of text-to-image generative models,” *arXiv preprint arXiv:2404.08030*, 2024.
- [68] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [69] “Pexels: The best free stock photos, royalty free images and videos shared by creators.” <https://www.pexels.com/>, 2023.
- [70] J. Jang, D. Yoon, S. Yang, S. Cha, M. Lee, L. Logeswaran, and M. Seo, “Knowledge unlearning for mitigating privacy risks in language models,” *arXiv preprint arXiv:2210.01504*, 2022.
- [71] J. Chen and D. Yang, “Unlearn what you want to forget: Efficient unlearning for llms,” *arXiv preprint arXiv:2310.20150*, 2023.
- [72] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. Iii, and K. Crawford, “Datasheets for datasets,” *Communications of the ACM*, vol. 64, no. 12, pp. 86–92, 2021.
- [73] K. Nichol, “Painter by numbers,” WikiArt, 2016, accessed: [insert date of access].
- [74] S. Liu, T. Lin, D. He, F. Li, M. Wang, X. Li, Z. Sun, Q. Li, and E. Ding, “Adaattn: Revisit attention mechanism in arbitrary neural style transfer,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 6649–6658.
- [75] X. Huang and S. Belongie, “Arbitrary style transfer in real-time with adaptive instance normalization,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1501–1510.
- [76] X. Li, S. Liu, J. Kautz, and M.-H. Yang, “Learning linear transformations for fast image and video style transfer,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3809–3817.

- [77] T. Q. Chen and M. Schmidt, “Fast patch-based style transfer of arbitrary style,” *arXiv preprint arXiv:1612.04337*, 2016.
- [78] D. Y. Park and K. H. Lee, “Arbitrary style transfer with style-attentional networks,” in *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 5880–5888.
- [79] J. An, S. Huang, Y. Song, D. Dou, W. Liu, and J. Luo, “Artflow: Unbiased image style transfer via reversible neural flows,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 862–871.
- [80] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [81] Q. Cai, M. Ma, C. Wang, and H. Li, “Image neural style transfer: A review,” *Computers and Electrical Engineering*, vol. 108, p. 108723, 2023.
- [82] A. Singh and P. Singh, “Image classification: a survey,” *Journal of Informatics Electrical and Electronics Engineering (JIEEE)*, vol. 1, no. 2, pp. 1–9, 2020.
- [83] Z. Zou, K. Chen, Z. Shi, Y. Guo, and J. Ye, “Object detection in 20 years: A survey,” *Proceedings of the IEEE*, 2023.
- [84] S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and D. Terzopoulos, “Image segmentation using deep learning: A survey,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 7, pp. 3523–3542, 2021.
- [85] L. Gatys, A. S. Ecker, and M. Bethge, “Texture synthesis using convolutional neural networks,” *Advances in neural information processing systems*, vol. 28, 2015.
- [86] L. A. Gatys, A. S. Ecker, and M. Bethge, “A neural algorithm of artistic style,” *arXiv preprint arXiv:1508.06576*, 2015.
- [87] Y. Deng, F. Tang, W. Dong, C. Ma, X. Pan, L. Wang, and C. Xu, “Stytr2: Image style transfer with transformers,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 11 326–11 336.
- [88] Y. Zhang, F. Tang, W. Dong, H. Huang, C. Ma, T.-Y. Lee, and C. Xu, “Domain enhanced arbitrary image style transfer via contrastive learning,” in *ACM SIGGRAPH 2022 Conference Proceedings*, 2022, pp. 1–8.
- [89] Y. Deng, F. Tang, W. Dong, H. Huang, C. Ma, and C. Xu, “Arbitrary video style transfer via multi-channel correlation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, 2021, pp. 1210–1217.
- [90] J. Huo, S. Jin, W. Li, J. Wu, Y.-K. Lai, Y. Shi, and Y. Gao, “Manifold alignment for semantically aligned style transfer,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14 861–14 869.
- [91] H. Chen, Z. Wang, H. Zhang, Z. Zuo, A. Li, W. Xing, D. Lu *et al.*, “Artistic style transfer with internal-external learning and contrastive learning,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 26 561–26 573, 2021.
- [92] D. Li, J. Li, and S. C. Hoi, “Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing,” *arXiv preprint arXiv:2305.14720*, 2023.
- [93] Z. Wang, Y. Jiang, Y. Lu, Y. Shen, P. He, W. Chen, Z. Wang, and M. Zhou, “In-context learning unlocked for diffusion models,” *arXiv preprint arXiv:2305.01115*, 2023.
- [94] X. Wang, W. Wang, Y. Cao, C. Shen, and T. Huang, “Images speak in images: A generalist painter for in-context visual learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 6830–6839.

- [95] Y. Bai, X. Geng, K. Mangalam, A. Bar, A. Yuille, T. Darrell, J. Malik, and A. A. Efros, “Sequential modeling enables scalable learning for large vision models,” *arXiv preprint arXiv:2312.00785*, 2023.
- [96] Q. Wang, X. Bai, H. Wang, Z. Qin, and A. Chen, “Instantid: Zero-shot identity-preserving generation in seconds,” *arXiv preprint arXiv:2401.07519*, 2024.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]** The introduction (Sec. 1) clearly specifies the scope of this paper and the major contributions.
  - (b) Did you describe the limitations of your work? **[Yes]** The limitations are discussed in Sec. 5
  - (c) Did you discuss any potential negative societal impacts of your work? **[No]** This work aims to improve the evaluation of machine unlearning technique, which will bring many positive societal impacts, such as mitigating copyright infringement, avoiding sexual or violent content generation, and mitigating the biases and stereotypes of the diffusion models. Given the technical focus of this work, there are no specific negative societal consequences directly stemming from it that need to be highlighted here.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? **[N/A]** There are no theoretical results in this paper.
  - (b) Did you include complete proofs of all theoretical results? **[N/A]** There are no theoretical results in this paper.
3. If you ran experiments (e.g. for benchmarks)...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** There are detailed instructions of the codes in the github repo to facilitate reproducing the results. The detailed experiment settings are also disclosed in Appx. B.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** Detailed training settings are disclosed in Sec. 4 and Appx. D.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[No]** The results reported in this work were averaged over dozens of individual unlearning cases, and the variance confirmed to be below 1%. We believe this indicates negligible variability. Given the large scale of the experiments, error bars were omitted to enhance readability and maintain focus on the key findings.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** The compute resources are disclosed in Appx. B.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? **[Yes]** The seed images used to curate UNLEARNCANVAS belong to the existing assets, and their sources are properly cited in Sec. 3.
  - (b) Did you mention the license of the assets? **[Yes]**
  - (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]** This paper includes a new dataset UNLEARNCANVAS and the link is released both in the abstract as well as in Appx. A.
  - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **[Yes]** The data being used is collected from an open-source photography website and this is mentioned in Sec. 3.



- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[Yes\]](#) The data does not contain any personally identifiable information or offensive content, as discussed in Appx. [A](#).
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#) No human subjects are involved.
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#) No human subjects are involved.
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#) No human subjects are involved.

# Appendix

## Table of Contents

<b>A UNLEARNCANVAS Dataset Details</b>	<b>19</b>
A.1 Motivation . . . . .	19
A.2 Composition: Styles and Object Classes . . . . .	19
A.3 Collection and Labeling Process . . . . .	19
A.4 Uses . . . . .	19
A.5 Distribution . . . . .	19
A.6 Maintenance . . . . .	19
A.7 Author Statements . . . . .	20
<b>B Reproducibility Statement and Detailed Experiment Settings</b>	<b>20</b>
B.1 Finetuning Style and Object Classifier with UNLEARNCANVAS . . . . .	20
B.2 Finetuning Stable Diffusion with UNLEARNCANVAS . . . . .	20
B.3 Implementation of DM Unlearning Methods Studied in This Work . . . . .	20
B.4 Evaluation Details of the Adversarial Prompt Generation for Unlearning Robustness	21
B.5 Experiment Details of the Style-Object Combination Unlearning . . . . .	22
B.6 Experiment Details of Sequential Unlearning . . . . .	23
B.7 Metrics Summary in All Unlearning Settings . . . . .	24
<b>C A Detailed Comparison between UNLEARNCANVAS and WIKIART</b>	<b>25</b>
<b>D Additional Experiment Results for Unlearning Evaluation</b>	<b>26</b>
D.1 Visualization of the Style and Object Unlearning Performance . . . . .	26
D.2 A Fine-Grained Comparison per Unlearning Target: ESD vs. SALUN . . . . .	26
D.3 Unlearning Heatmaps of More Unlearning Methods . . . . .	26
<b>E Visualizations</b>	<b>29</b>
E.1 Illustrations of the Styles and Objects in UNLEARNCANVAS . . . . .	29
E.2 Visualization of Style Unlearning . . . . .	30
E.3 Visualization of the Unlearning Performance in the Presence of Adversarial Prompts	30
<b>F Broader Use Cases of UNLEARNCANVAS</b>	<b>34</b>
F.1 Benchmarking Style Transfer using UNLEARNCANVAS . . . . .	34
F.2 Other Possible Applications of UNLEARNCANVAS . . . . .	36
<b>G Impact Statement</b>	<b>37</b>

## A UNLEARNCANVAS Dataset Details

In this section, we provide a detailed description of the dataset following ‘datasheet for datasets’ [72].

### A.1 Motivation

This dataset is collected to enable a precise, comprehensive, and automated quantitative evaluation framework for MU (machine unlearning) methods in DMs (diffusion models). The current evaluation plans used in the existing literature have exposed several weaknesses, which may lead to incomplete, inaccurate or even biased results, see a more detailed discussion in Sec. 2. To the best of our knowledge, there are no datasets specifically designed to meet the assessing requirements of DM unlearning. Therefore, UNLEARNCANVAS is designed, collected, and made to fill in this gap.

### A.2 Composition: Styles and Object Classes

There are 60 predetermined artistic styles provided by Fotor [11]. The images in the same artistic style all share high stylistic consistency, which enables a high-precision style classifier to be trained on them. In Fig. A9, we list some examples of the images in each style to illustrate these styles. There are 20 distinct object classes in UNLEARNCANVAS. In Fig. A10, we list some examples of the images in each object to illustrate these classes.

### A.3 Collection and Labeling Process

The construction of UNLEARNCANVAS involves two main steps: seed image collection and subsequent image stylization; see Fig. 3 for an illustration. For seed image collection, a set of high-resolution real-world photos are collected from Pexels [69], providing open-sourced photographs. There are 20 seed images collected for each of the 20 object classes; see Fig. A10. After collecting the seed images, we stylize each and every seed image into all 60 predetermined artistic styles; see Fig. A9 with Fotor [11]. After the stylization of all the images, the dataset is structured in a hierarchical manner, and each image is labeled with both its style and object classes. In order to support text-to-image training, each image is annotated with the prompt ‘An image of *object* in *style*.’.

### A.4 Uses

UNLEARNCANVAS can be used to evaluate MU methods in different unlearning scenarios. Please see Sec. 4 for more details. In addition, we stress that UNLEARNCANVAS can be used for more real-world tasks than unlearning, and we provide an example of how UNLEARNCANVAS can be used to systematically evaluate another generative task, style transfer, in Appx. F. We provide very detailed instructions with codes in the GitHub code repository.

### A.5 Distribution

UNLEARNCANVAS is an open-sourced dataset and is based on the existing open-sourced data [69]. We make access to the dataset public under the MIT license. We remark that no personally identifiable information or offensive content is included in this dataset. The dataset can be accessed either through Google Drive and HuggingFace. More resources on the dataset, such as the introduction video and the benchmark, can be found in the official *project webpage*.

### A.6 Maintenance

The dataset will be maintained by the lead author Yihua Zhang. If needed, the email for contacting is zhan1908@msu.edu. The dataset may be updated if needed (with the inclusion of more seed images, more artistic styles, and more objects). The updates will be ad-hoc and will not be periodical. Each time the dataset is updated, the updates will be reflected in the same GitHub code repository.

## A.7 Author Statements

The collector and the lead author of this dataset, Yihua Zhang, bears full responsibility for any violation of rights that may arise from the collection of the data included in this research.

## B Reproducibility Statement and Detailed Experiment Settings

In this section, we provide detailed instructions on the reproduction of our results in Sec. 4, including the settings of training, the implementation details of the tested machine unlearning methods, and the evaluation details in each unlearning scenario. We denote in Sec. 4, 50 out of the 60 styles proposed in UNLEARNCANVAS are studied and evaluated, with the other 10 styles being intentionally reserved for future explorations, such as generalizations.

### B.1 Finetuning Style and Object Classifier with UNLEARNCANVAS

Style and object classifiers need to be trained as part of the testbed proposed in our evaluation pipeline (Fig. 5). Here, we adopted a ViT-L/16 model [68] pretrained on ImageNet and finetune it on UNLEARNCANVAS. UNLEARNCANVAS are split into the train set and test set with a ratio of 9 : 1. After hyper-parameter tuning, the classifiers are trained with Adam optimizer at a learning rate of 0.01 for 10 epochs.

### B.2 Finetuning Stable Diffusion with UNLEARNCANVAS

The other part of the testbed is a diffusion model capable of generating high quality images in all the styles associated with all the objects encompassed in UNLEARNCANVAS in order to guarantee a trustworthy and unbiased evaluation.

**Training settings.** Practically, we finetune the pretrained Stable Diffusion (SD) v1.5 on UNLEARNCANVAS for 20k steps with a learning rate of  $1e - 6$ . Unless otherwise stated, we strictly follow the training configurations used in Stable Diffusion [1]. For each image in UNLEARNCANVAS, we annotate the data with text prompt An image of  $\{object\}$  in  $\{style\}$ , where the  $object$  and  $style$  are the corresponding object and style label. For seed images, the  $style$  label we use is ‘photo’ style. We use the training scripts provided by Diffuser official tutorial (<https://huggingface.co/docs/diffusers/v0.13.0/en/training/text2image>) and the pretraining model card is [runwayml/stable-diffusion-v1-5](https://huggingface.co/runwayml/stable-diffusion-v1-5). During training, the checkpoints will be saved every 1000 steps.

**Evaluation.** To evaluate the quality of the saved checkpoints and select the best one for unlearning study, the checkpoints are first used to generate an image set with the same prompt as training (An image of  $\{object\}$  in  $\{style\}$ ) by traversing all the possible style and object labels. Each prompt are used to generate 5 images with different random seed. Each image are sampled with 100 steps with a guidance coefficient of 9. The image set for each checkpoint are fed into the style and object classifier trained in Appx. B.1. The model with the highest average performance on all the styles and objects are selected as the testbed for MU study. The classification performance are also used as a reference for later IRA/CRA comparison, which are disclosed in the first row of Fig. 6 (left).

**Computing resource.** In this work, we employ  $40 \times$  NVIDIA RTX A6000 GPUs to conduct all the model training, unlearning, image generation, and evaluations. When we finetuned the StableDiffusion on UNLEARNCANVAS,  $8 \times$  GPUs were used for parallel computing. Other experiments were all carried out in a single-GPU environment. Around 60,000 GPU hours in total were spent to complete all the experiments.

### B.3 Implementation of DM Unlearning Methods Studied in This Work

In this work, we inspected a series of stateful MU methods for DMs. For each method, we use their publicly released source codes as code bases, which are listed below:

- ESD [23]: <https://github.com/rohitgandikota/erasing>

- CA [25]: <https://github.com/nupurkmr9/concept-ablation>
- UCE [24]: <https://github.com/rohitgandikota/unified-concept-editing>
- FMN [28]: <https://github.com/SHI-Labs/Forget-Me-Not>
- SalUn: [27]: <https://github.com/OPTML-Group/Unlearn-Saliency>
- SEOT: [30]: <https://github.com/sen-mao/SuppressEOT>
- SPM: [26]: <https://github.com/Con6924/SPM>
- EDiff: [31]: <https://github.com/JingWu321/EraseDiff>
- SHS: [32]: <https://github.com/JingWu321/Scissorhands>

In particular, we adopt the following training settings to adapt the methods to our dataset:

- ESD: Based on the suggestions from the authors, ESD is used to only finetune the cross attention-related model weights (ESD-x). Other settings strict follow the ones used in the paper.
- CA: In order to ablate concepts using CA, we first use ChatGPT to generate a list of simple prompts for each concept (including the styles and the objects), namely anchor prompts. Each anchor prompt is a simple one-sentence description of the unlearning target.
- UCE: This method requires a guided concept (prompt) for each unlearning concept. For style unlearning, we use the prompt An image in  $\{style^*\}$  as the guided concept, where  $style^*$  represents the next style in UNLEARNCANVAS in alphabetical order. Similarly, An image of  $object^*$  is used for object unlearning, where  $object^*$  is the next object in UNLEARNCANVAS in alphabetical order.
- FMN: This method requires the images associated with the unlearning target. For simplicity and best performance, we randomly select 20 images associated with the unlearning concept. For the first stage of FMN, we run text inversion for 500 steps with a learning rate of  $1e-4$ , and for the second step, we used the inversed text to unlearn the cross attention layers of the model for 100 steps. The hyper-parameters of learning rate, maximum steps, and tunable parameters (cross-attention or non-cross-attention) are carefully tuned with grid search.
- SalUn: This method involves two steps, the mask finding (weight saliency analysis) and the model unlearning. For mask finding, we tuned the mask ratio, while for unlearning, we tune the hyper-parameter learning rate and unlearning intensity. All the parameters are tuned with grid search. For both steps, the mask or model is trained with 10 epochs.
- SEOT: To generate unlearned images, we use the prompt An  $\{object^*\}$  image in  $\{style^*\}$ . We then suppress either  $object^*$  or  $style^*$  individually. Other settings strict follow the ones used in the paper.
- SPM: Following the hyperparameters provided by the authors, we trained and obtained Pre-tuned SPMs for all  $object^*$  and  $style^*$ . During image generation, we combine the pre-tuned SPMs with the DM. By calculating the association between words in the prompt and the target word, we determine whether to allow the specified word to preserve, and generate the corresponding image.
- EDiff: Based on the authors’ suggestions, EDiff is used to finetune only the cross-attention-related model weights (EraseDiff-x). During the unlearning process, we adjusted the hyperparameters, specifically the learning rate and the number of unlearning epochs. The model is trained with 5 epochs.
- SHS: SHS consists of two stages: trimming and repairing. During the trimming stage, certain weights are re-initialized. The repairing stage then restores the model’s utility. Throughout the unlearning process, we finetuned the hyperparameters, focusing on the learning rate and the number of unlearning epochs. Ultimately, we selected 2 epochs as the optimal number.

#### B.4 Evaluation Details of the Adversarial Prompt Generation for Unlearning Robustness

In Sec. 4, we evaluated the robustness of different MU methods against adversarial prompts. Here, we use the state-of-the-art method, UnlearnDiffAtk [59] to generate adversarial prompts. We set the prepended prompt perturbations by  $N = 5$  tokens for both style and object unlearning. Following the

original attack setting in UnlearnDiffAtk [59], to optimize the adversarial perturbations, we sample 50 diffusion time steps and perform PGD running for 40 iterations with a learning rate of 0.01 at each step. Prior to projection onto the discrete text space, we utilize the AdamW optimizer.

## B.5 Experiment Details of the Style-Object Combination Unlearning

**Unlearning targets.** In Sec. 4, we evaluate the capability of different MU methods on performing unlearning at a finer scale, and we use the style-object combinations as unlearning targets for evaluation. Ideally, the UNLEARNCANVAS dataset can generate 1200 ( $60 \times 20$ ) style-object combinations. In this work, we randomly select 50 of these combinations for evaluation and we list these combinations below. For each method, the same hyper-parameters are used for each MU method as the ones for style and object unlearning in Tab. 2. The unlearning targets include:

- ‘An image of Architectures in Abstractionism style.’
- ‘An image of Bears in Artist Sketch style.’
- ‘An image of Birds in Blossom Season style.’
- ‘An image of Butterfly in Bricks style.’
- ‘An image of Cats in Byzantine style.’
- ‘An image of Dogs in Cartoon style.’
- ‘An image of Fishes in Cold Warm style.’
- ‘An image of Flame in Color Fantasy style.’
- ‘An image of Flowers in Comic Etch style.’
- ‘An image of Frogs in Crayon style.’
- ‘An image of Horses in Cubism style.’
- ‘An image of Human in Dadaism style.’
- ‘An image of Jellyfish in Dapple style.’
- ‘An image of Rabbits in Defoliation style.’
- ‘An image of Sandwiches in Early Autumn style.’
- ‘An image of Sea in Expressionism style.’
- ‘An image of Statues in Fauvism style.’
- ‘An image of Towers in French style.’
- ‘An image of Trees in Glowing Sunset style.’
- ‘An image of Waterfalls in Gorgeous Love style.’
- ‘An image of Architectures in Greenfield style.’
- ‘An image of Bears in Impressionism style.’
- ‘An image of Birds in Ink Art style.’
- ‘An image of Butterfly in Joy style.’
- ‘An image of Cats in Liquid Dreams style.’
- ‘An image of Dogs in Magic Cube style.’
- ‘An image of Fishes in Meta Physics style.’
- ‘An image of Flame in Meteor Shower style.’
- ‘An image of Flowers in Monet style.’
- ‘An image of Frogs in Mosaic style.’
- ‘An image of Horses in Neon Lines style.’
- ‘An image of Human in On Fire style.’
- ‘An image of Jellyfish in Pastel style.’
- ‘An image of Rabbits in Pencil Drawing style.’

- ‘An image of Sandwiches in Picasso style.’
- ‘An image of Sea in Pop Art style.’
- ‘An image of Statues in Red Blue Ink style.’
- ‘An image of Towers in Rust style.’
- ‘An image of Waterfalls in Sketch style.’
- ‘An image of Architectures in Sponge Dabbed style.’
- ‘An image of Bears in Structuralism style.’
- ‘An image of Birds in Superstring style.’
- ‘An image of Butterfly in Surrealism style.’
- ‘An image of Cats in Ukiyoe style.’
- ‘An image of Dogs in Van Gogh style.’
- ‘An image of Fishes in Vibrant Flow style.’
- ‘An image of Flame in Warm Love style.’
- ‘An image of Flowers in Warm Smear style.’
- ‘An image of Frogs in Watercolor style.’
- ‘An image of Horses in Winter style.’

**Evaluation.** The evaluation of the style-object combination unlearning concerns four quantitative metrics, one for unlearning effectiveness and three for retainability. Before the evaluation, an answer set will be generated exactly following the same procedure introduced in Sec. 3 and Fig. 5 after unlearning each target. First, the UA (unlearning accuracy) will be evaluated for each answer set, which stands for the ratio of images generated by the target prompt that are neither classified into the target object nor the target style class. A high UA denotes a better ability to successfully unlearn the target combination. Second, the retainability of generation associated with those prompts close to the unlearning target prompt will be evaluated. These prompts can be divided into two groups, the ones sharing the same style but not the object class and the ones sharing the object but not the style class. The classification accuracy of the former corresponds to the retainability of the style, *i.e.*, style consistency (SC), while the latter one denotes the object consistency (OC). These two quantitative metrics evaluate how well the unlearning method precisely define the unlearning scope and retain the generation ability of those close but innocent concept. Thirdly, the retainability of the rest unrelated prompts (UP) are evaluated, which is the last quantitative evaluation metric. The results reported in Tab. 3 are averaged over all the unlearning cases shown above.

## B.6 Experiment Details of Sequential Unlearning

In Sec. 4, we also evaluated the MU methods with the task of sequential unlearning (SU), where the efficacy of MU methods in handling multiple sequential unlearning requests  $\{\mathcal{T}_i\}$  are evaluated. This requires models not only to unlearn new targets effectively but also to maintain the unlearning of previous targets, while retaining all other knowledge. In the experiments, 6 styles are randomly selected as the unlearning targets and excluded from the RA evaluation. The UA of all the already unlearned target will be assessed each time a new request is accomplished. The selected 6 styles include:

- Abstractionism
- Byzantine
- Cartoon
- Cold Warm
- Ukiyoe
- Van Gogh

After each unlearning request, the unlearning effectiveness and retainability are evaluated. Specifically, the unlearning accuracy of all the unlearning targets in the previous requests are evaluated to evaluate how the unlearning effect lasts when new unlearning requests arrive. In the meantime, the retainability of all the other concepts that are not selected as unlearning targets are evaluated, and to ease the presentation, the retain accuracy of all the concepts (styles and objects) are averaged and reported.

### B.7 Metrics Summary in All Unlearning Settings

Besides the UNLEARNCANVAS dataset, the various quantitative evaluation metrics proposed in this work are part of the major contributions to a comprehensive and precise evaluation for DM unlearning methods. As there are various unlearning scenarios studied in this work, we provide a summary of these metrics in Tab. A1, including their abbreviations, descriptions, and related tables or figures.

Table A1: A summary of the quantitative metrics used in this work, including their abbreviations, meanings and where they are used.

Metrics	Description	Usages (Table & Figure)
Style/Object Unlearning		
UA	Unlearning accuracy	Fig. 1, Tab. 1
IRA	In-domain unlearning accuracy	Fig. 1, Tab. 1
CRA	Cross-domain unlearning accuracy	Fig. 1, Tab. 1
Unlearning Robustness against Adversarial Prompts		
Rob.	Unlearning robustness, unlearning accuracy in the presence of adversarial prompts	Fig. 1
Style-Object Combination Finer-Scale Unlearning		
FU/UA	Unlearning accuracy in finer-scale unlearning	Fig. 1, Tab. 3
SC	Retainability evaluation of style consistency	Tab. 3
OC	Retainability evaluation of object consistency	Tab. 3
UP	Retainability evaluation of unrelated prompts	Tab. 3
FR	Retainability evaluation in finer-scale unlearning, averaged by SC, OC, and UP	Fig. 1
Sequential or Continual Unlearning		
SU or CU	Unlearning accuracy in the context of sequential unlearning	Fig. 1, Tab. A5
SR or CR	Retainability in the context of sequential unlearning	Fig. 1, Tab. A5



## C A Detailed Comparison between UNLEARNCANVAS and WIKIART



Figure A1: Image examples with the same style label from WIKIART [66] and UNLEARNCANVAS. Images of the same artistic style in UNLEARNCANVAS exhibit high stylistic consistency compared to WIKIART.

**UNLEARNCANVAS vs. WIKIART.** To the best of our knowledge, WIKIART [73] is the most relevant baseline dataset to ours. In **Tab. A2**, we provide a direct comparison of the key attributes of these two datasets. UNLEARNCANVAS differs from WIKIART in the following aspects.

**First**, UNLEARNCANVAS includes a greater number of high-resolution images (15M) compared to WIKIART (2M), a factor that may enhance the training of state-of-the-art DMs.

Table A2: Comparison with WIKIART, the most relevant dataset containing stylized images to ours. UNLEARNCANVAS stands out notably from WIKIART due to its characteristics of being supervised, balanced, and maintaining high stylistic cohesiveness.

Dataset	Resolution (Pixels/Image)	Style-wise Supervised	High Stylistic Consistency	Class-wise Balanced
WIKIART [73]	~ 2M	✗	✗	Style-wise ✗ Object-wise ✗
UNLEARNCANVAS	~ 15M	✓	✓	Style-wise ✓ Object-wise ✓

**Second**, UNLEARNCANVAS surpasses WIKIART in terms of both intra-style coherence and inter-style distinctiveness, as illustrated in **Fig. A1**, where images labeled with ‘Van Gogh Style’ from both datasets are compared. In UNLEARNCANVAS, the images exhibit high stylistic consistency, while WIKIART lacks the necessary clarity for precise assessment. This will hamper the MU evaluation as discussed in the challenges (C2) and (C3). This benefits can also be reflected by the training performance using UNLEARNCANVAS and WIKIART. The results are reported in **Tab. A3** and **Tab. A4**, respectively. As we can see, the classifier is much more easily trained on UNLEARNCANVAS, justifying the higher discernible features within each style in UNLEARNCANVAS.

Table A3: Art style reproduction quality using SD v1.5 and SD v2.0 finetuned on WIKIART. Images are generated with the prompt “A painting in *artist* style”, where *artist* refers to those included in WIKIART. The test accuracy on DM-generated images and original WIKIART test images is reported using the style classifier finetuned from the pretrained ViT-L/16 on WIKIART.

Image Source	Images by SD v1.5	Images by SD v2.0	WIKIART Test Set
Accuracy	41.2%	56.7%	85.4%

**Third**, the images in UNLEARNCANVAS are style-wise supervised. For each seed image, a stylized counterpart can be found in each style class. This is beneficial for tasks other than unlearning for text-to-image task, such as image editing, image stylization, and style transfer, which can provide a ground truth image for precise and robust evaluation. This will be detailed in **Appx. F**.

Table A4: Style classification results of a ViT-Large [68] as a style classifier trained on UNLEARN-CANVAS. After convergence, the classifier is tested on the test set and the image set generated by SD v1.5 finetuned on UNLEARNCANVAS.

	UNLEARNCANVAS Train Set	UNLEARNCANVAS Test Set	Images by SD v1.5 tuned on UNLEARNCANVAS
Accuracy	100.0%	99.9%	98.8%

## D Additional Experiment Results for Unlearning Evaluation

### D.1 Visualization of the Style and Object Unlearning Performance

To make a more direct comparison among different MU methods reported in Tab. 1, the results are visualized in the radar chart Fig. A2. This figure illustrates that no method dominates across all assessment dimensions. This underscores the complexity of unlearning in generative models and the need for further improvement.

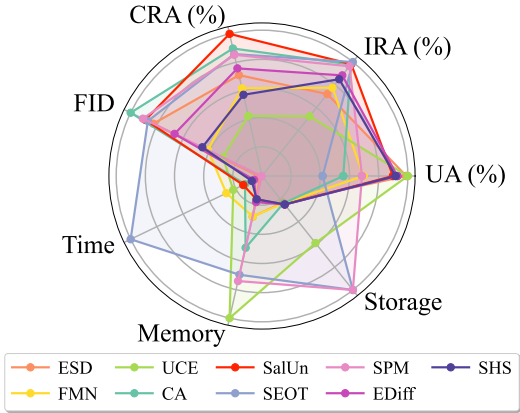


Figure A2: Performance visualization for various unlearning methods as summarized in Table 2. For UA, IRA, and CRA, the results are averaged over the style and object unlearning scenarios. Other metrics undertake the inverse operation as a smaller values represent better performance. Results are normalized to 0% ~ 100% per metric.

### D.2 A Fine-Grained Comparison per Unlearning Target: ESD vs. SALUN

Following the analysis of ESD in Fig. 6, we next turn our focus to a comparative analysis with SALUN, a method that demonstrated a better balance between unlearning and retaining according to Tab. 2. A similar accuracy heatmap for SALUN is presented in Fig. A3. Compared to ESD, SALUN exhibits more consistent performance across various unlearning scenarios, as indicated by the more uniform color distribution in the heatmap. This also suggests enhanced retainability. However, it is noticeable that SALUN does not reach the same level of UA (Unlearning Accuracy) as ESD, as evidenced by the darker diagonal values in Fig. A3. This observation reinforces the existence of a trade-off between unlearning effectiveness and retainability in the visual generative MU task, a phenomenon paralleled in other tasks such as classification.

### D.3 Unlearning Heatmaps of More Unlearning Methods

In Fig. A4~Fig. A8, we provide more unlearning heatmap visualizations in the same format as Fig. 6 and Fig. A3 in order to provide a more detailed unlearning performance dissection for all the DM unlearning methods studied in this work.

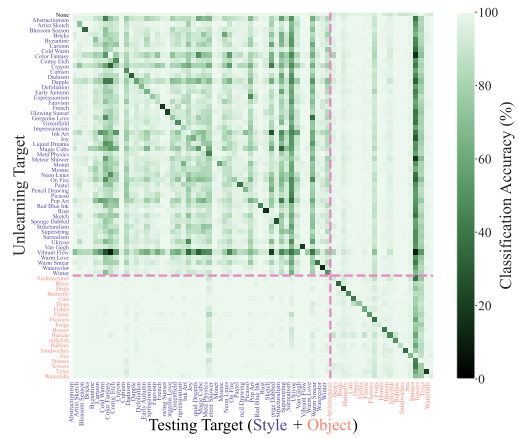


Figure A3: Heatmap visualization of SalUn. The plot setting is identical to Figure 6.

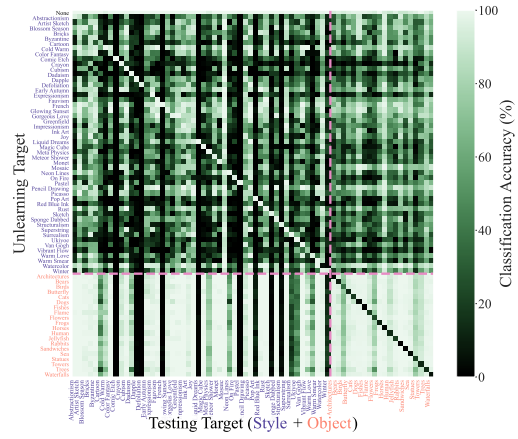


Figure A4: Heatmap visualization of FMN. The plot setting is identical to Figure 6.

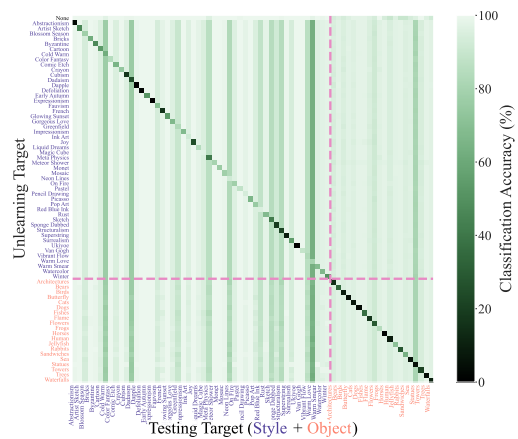


Figure A5: Heatmap visualization of SEOT. The plot setting is identical to Figure 6.

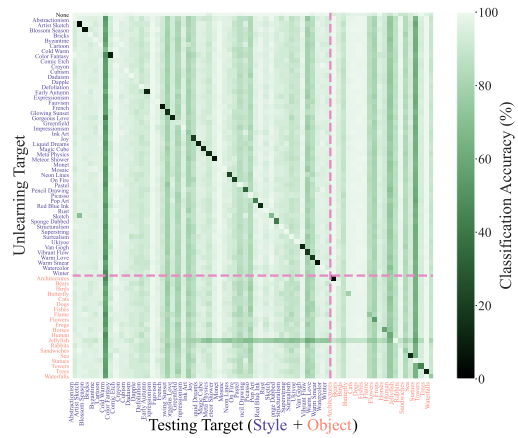


Figure A6: Heatmap visualization of SPM. The plot setting is identical to Figure 6.

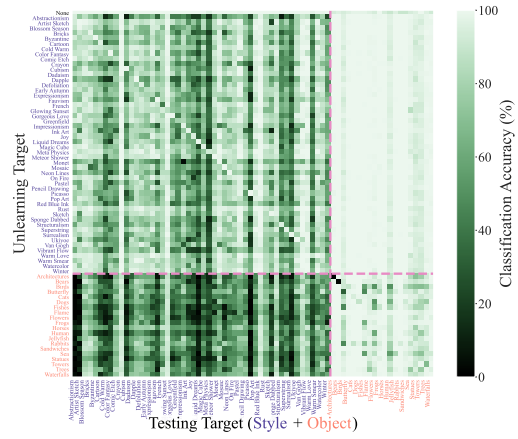


Figure A7: Heatmap visualization of Ediff. The plot setting is identical to Figure 6.

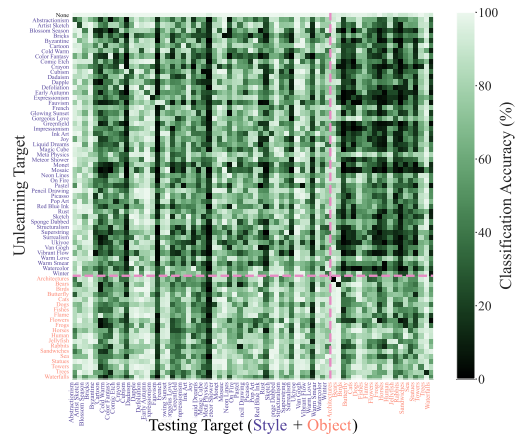


Figure A8: Heatmap visualization of SHS. The plot setting is identical to Figure 6.

Table A5: Performance comparison of different DM unlearning methods in the sequential unlearning setting. Each column represents a new unlearning request, denoted by  $\mathcal{T}_i$ , where  $\mathcal{T}_1$  is the oldest. Each row represents the UA for a specific unlearning objective or the retaining accuracy (RA), given by the average of IRA and CRA. Results indicating *unlearning rebound* effect are highlighted in orange, and those signifying *catastrophic retaining failure* are marked in red.

Method: ESD							Method: FMN							
Metrics	$\mathcal{T}_1$	$\mathcal{T}_1 \sim \mathcal{T}_2$	$\mathcal{T}_1 \sim \mathcal{T}_3$	$\mathcal{T}_1 \sim \mathcal{T}_4$	$\mathcal{T}_1 \sim \mathcal{T}_5$	$\mathcal{T}_1 \sim \mathcal{T}_6$	Metrics	$\mathcal{T}_1$	$\mathcal{T}_1 \sim \mathcal{T}_2$	$\mathcal{T}_1 \sim \mathcal{T}_3$	$\mathcal{T}_1 \sim \mathcal{T}_4$	$\mathcal{T}_1 \sim \mathcal{T}_5$	$\mathcal{T}_1 \sim \mathcal{T}_6$	
	Unlearning Request							Unlearning Request						
UA	$\mathcal{T}_1$	100%	99%	95%	87%	81%	75%	$\mathcal{T}_1$	88%	99%	99%	98%	99%	99%
	$\mathcal{T}_2$	-	100%	100%	96%	87%	79%	$\mathcal{T}_2$	-	95%	99%	99%	98%	99%
	$\mathcal{T}_3$	-	-	100%	98%	99%	98%	$\mathcal{T}_3$	-	-	97%	98%	99%	99%
	$\mathcal{T}_4$	-	-	-	100%	100%	99%	$\mathcal{T}_4$	-	-	-	99%	99%	99%
	$\mathcal{T}_5$	-	-	-	-	100%	99%	$\mathcal{T}_5$	-	-	-	-	99%	99%
	$\mathcal{T}_6$	-	-	-	-	-	100%	$\mathcal{T}_6$	-	-	-	-	-	100%
RA	77.46%	52.94%	35.99%	24.86%	18.69%	12.95%	RA	82.39%	14.56%	13.34%	10.42%	9.83%	8.76%	
Method: UCE							Method: CA							
Metrics	$\mathcal{T}_1$	$\mathcal{T}_1 \sim \mathcal{T}_2$	$\mathcal{T}_1 \sim \mathcal{T}_3$	$\mathcal{T}_1 \sim \mathcal{T}_4$	$\mathcal{T}_1 \sim \mathcal{T}_5$	$\mathcal{T}_1 \sim \mathcal{T}_6$	Metrics	$\mathcal{T}_1$	$\mathcal{T}_1 \sim \mathcal{T}_2$	$\mathcal{T}_1 \sim \mathcal{T}_3$	$\mathcal{T}_1 \sim \mathcal{T}_4$	$\mathcal{T}_1 \sim \mathcal{T}_5$	$\mathcal{T}_1 \sim \mathcal{T}_6$	
	Unlearning Request							Unlearning Request						
UA	$\mathcal{T}_1$	93%	95%	98%	96%	97%	98%	$\mathcal{T}_1$	58%	55%	59%	45%	44%	40%
	$\mathcal{T}_2$	-	97%	98%	98%	98%	95%	$\mathcal{T}_2$	-	76%	58%	51%	47%	44%
	$\mathcal{T}_3$	-	-	95%	97%	98%	99%	$\mathcal{T}_3$	-	-	45%	41%	40%	37%
	$\mathcal{T}_4$	-	-	-	98%	98%	98%	$\mathcal{T}_4$	-	-	-	71%	70%	60%
	$\mathcal{T}_5$	-	-	-	-	97%	99%	$\mathcal{T}_5$	-	-	-	-	69%	51%
	$\mathcal{T}_6$	-	-	-	-	-	99%	$\mathcal{T}_6$	-	-	-	-	-	57%
RA	81.42%	29.38%	18.72%	15.34%	13.32%	11.31%	RA	97.24%	93.39%	84.46%	79.32%	71.40%	60.53%	
Method: SalUn							Method: SPM							
Metrics	$\mathcal{T}_1$	$\mathcal{T}_1 \sim \mathcal{T}_2$	$\mathcal{T}_1 \sim \mathcal{T}_3$	$\mathcal{T}_1 \sim \mathcal{T}_4$	$\mathcal{T}_1 \sim \mathcal{T}_5$	$\mathcal{T}_1 \sim \mathcal{T}_6$	Metrics	$\mathcal{T}_1$	$\mathcal{T}_1 \sim \mathcal{T}_2$	$\mathcal{T}_1 \sim \mathcal{T}_3$	$\mathcal{T}_1 \sim \mathcal{T}_4$	$\mathcal{T}_1 \sim \mathcal{T}_5$	$\mathcal{T}_1 \sim \mathcal{T}_6$	
	Unlearning Request							Unlearning Request						
UA	$\mathcal{T}_1$	84%	79%	78%	65%	67%	64%	$\mathcal{T}_1$	55%	59%	50%	49%	47%	48%
	$\mathcal{T}_2$	-	81.42%	75%	72%	69%	61%	$\mathcal{T}_2$	-	62%	59%	58%	60%	63%
	$\mathcal{T}_3$	-	-	90%	85%	84%	87%	$\mathcal{T}_3$	-	-	42%	39%	40%	41%
	$\mathcal{T}_4$	-	-	-	84%	86%	81%	$\mathcal{T}_4$	-	-	-	57%	59%	60%
	$\mathcal{T}_5$	-	-	-	-	79%	81%	$\mathcal{T}_5$	-	-	-	-	51%	51%
	$\mathcal{T}_6$	-	-	-	-	-	89%	$\mathcal{T}_6$	-	-	-	-	-	43%
RA	85.43%	80.32%	71.42%	65.41%	63.24%	60.19%	RA	72.39%	70.42%	67.89%	60.45%	55.32%	51.12%	
Method: EDiff							Method: SHS							
Metrics	$\mathcal{T}_1$	$\mathcal{T}_1 \sim \mathcal{T}_2$	$\mathcal{T}_1 \sim \mathcal{T}_3$	$\mathcal{T}_1 \sim \mathcal{T}_4$	$\mathcal{T}_1 \sim \mathcal{T}_5$	$\mathcal{T}_1 \sim \mathcal{T}_6$	Metrics	$\mathcal{T}_1$	$\mathcal{T}_1 \sim \mathcal{T}_2$	$\mathcal{T}_1 \sim \mathcal{T}_3$	$\mathcal{T}_1 \sim \mathcal{T}_4$	$\mathcal{T}_1 \sim \mathcal{T}_5$	$\mathcal{T}_1 \sim \mathcal{T}_6$	
	Unlearning Request							Unlearning Request						
UA	$\mathcal{T}_1$	97%	93%	91%	93%	85%	90%	$\mathcal{T}_1$	81%	73%	74%	93%	94%	97%
	$\mathcal{T}_2$	-	92%	89%	93%	91%	87%	$\mathcal{T}_2$	-	69%	61%	89%	94%	97%
	$\mathcal{T}_3$	-	-	96%	93%	90%	84%	$\mathcal{T}_3$	-	-	75%	92%	96%	90%
	$\mathcal{T}_4$	-	-	-	91%	92%	90.22%	$\mathcal{T}_4$	-	-	-	91%	95%	97%
	$\mathcal{T}_5$	-	-	-	-	99%	97%	$\mathcal{T}_5$	-	-	-	-	92%	96%
	$\mathcal{T}_6$	-	-	-	-	-	94%	$\mathcal{T}_6$	-	-	-	-	-	94%
RA	92.34%	89.37%	14.35%	12.31%	12.82%	7.42%	RA	88.41%	84.32%	73.98%	69.19%	10.76%	10.11%	

## E Visualizations

In this section, we aim to provide plenty of visualizations, illustrations, and qualitative results of the dataset and the quantitative results shown in Sec. 4 and Appx. D. These visualizations are intended to deepen the understanding of the effects of different MU methods and clearly illustrate the challenges identified in previous sections. We hope these visual aids will enable readers to more effectively grasp the nuances of MU methods and their implications for DMs (Diffusion Models).

### E.1 Illustrations of the Styles and Objects in UNLEARNCANVAS

We first provide an illustration of the styles and object classes included in UNLEARNCANVAS. Specifically, we show the styles in Fig. A9 and object classes in Fig. A10 with the style and object names disclosed in the captions.



Figure A9: An illustration of the images in each style in UNLEARNCANVAS used in, which are stylized from the same seed image from the ‘Dogs’ object class. The seed image is presented in Fig. A10 (6). Images are cropped and down-scaled for illustration purpose. The name of the styles are: (1) Abstractionism; (2) Artist Sketch; (3) Blossom Season; (4) Blue Blooming; (5) Bricks; (6) Byzantine; (7) Cartoon; (8) Cold Warm; (9) Color Fantasy; (10) Comic Etch; (11) Crayon; (12) Crypto Punks; (13) Cubism; (14) Dadaism; (15) Dapple; (16) Defoliation; (17) Dreamweave; (18) Early Autumn; (19) Expressionism; (20) Fauvism; (21) Foliage Patchwork; (22) French; (23) Glowing Sunset; (24) Gorgeous Love; (25) Greenfield; (26) Impasto; (27) Impressionism; (28) Ink Art; (29) Joy; (30) Liquid Dreams; (31) Palette Knife; (32) Magic Cube; (33) Meta Physics; (34) Meteor Shower; (35) Monet; (36) Mosaic; (37) Neon Lines; (38) On Fire; (39) Pastel; (40) Pencil Drawing; (41) Picasso; (42) Pointillism; (43) Pop Art; (44) Rainwash; (45) Realistic Watercolor; (46) Red Blue Ink; (47) Rust; (48) Sketch; (49) Sponge Dabbed; (50) Structuralism; (51) Superstring; (52) Surrealism; (53) Techno; (54) Ukiyoe; (55) Van Gogh; (56) Vibrant Flow; (57) Warm Love; (58) Warm Smear; (59) Watercolor; (60) Winter.

## E.2 Visualization of Style Unlearning

In Fig. A11, we provide abundant generation examples of all the 9 methods benchmarked in this work in a case study of unlearning the ‘Cartoon’ style. Both the successful and failure cases are demonstrated in the context of unlearning effectiveness, in-domain retainability, and cross-domain retainability.

## E.3 Visualization of the Unlearning Performance in the Presence of Adversarial Prompts

In Fig. A12, we provide visualizations for the effect of adversarial prompts. As revealed in Fig. 8, all the DM unlearning methods experience a significant drop in unlearning effectiveness when attacked by the adversarial prompt, enabled by UnlearnDiffAtk [59]. We provide the image generation in four unlearning cases (two for style unlearning and two for object unlearning), and show the images of the unlearning target successfully generated in the presence of adversarial prompts.

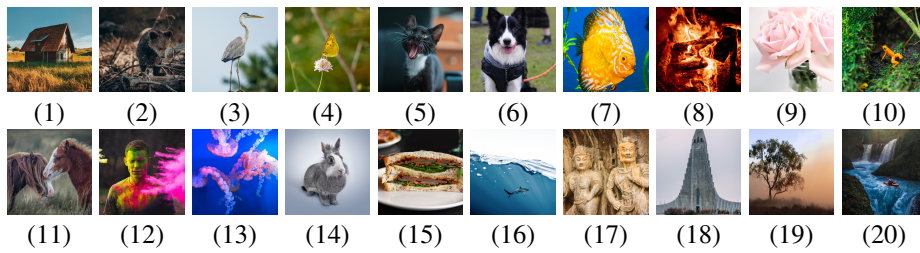
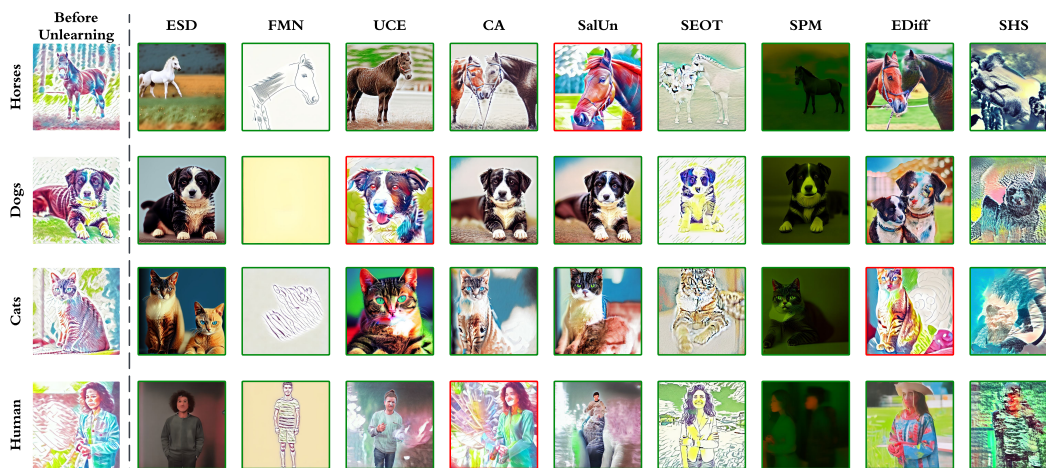


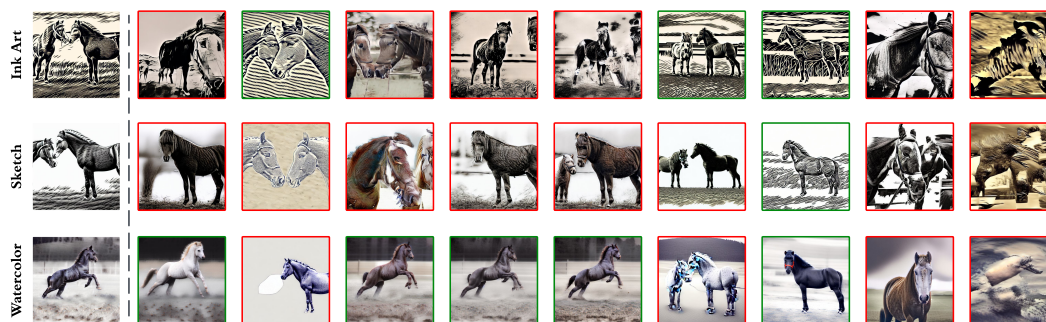
Figure A10: An illustration of the seed images in each object class in UNLEARNCANVAS. Images are cropped and down-scaled for illustration purpose. The name of the object classes are: (1) Architecture; (2) Bear; (3) Bird; (4) Butterfly; (5) Cat; (6) Dog; (7) Fish; (8) Flame; (9) Flowers; (10) Frog; (11) Horse; (12) Human; (13) Jellyfish; (14) Rabbits; (15) Sandwich; (16) Sea; (17) Statue; (18) Tower; (19) Tree; (20) Waterfalls.

### Unlearning Target Concept - **Cartoon** style

Unlearning Effectiveness Evaluation: Test Prompt Template: "An image of {object} in **Cartoon** style"



In-Domain Retainability Evaluation: Test Prompt Template: "An image of Horses in [style] style"



Cross-Domain Retainability Evaluation: Test Prompt Template: "An image of {object}."

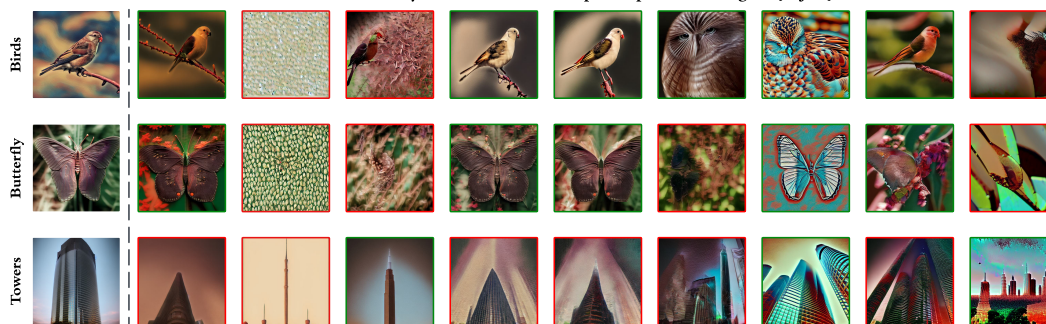


Figure A11: visualization of the unlearning performance of different methods on the task of style unlearning. Three text prompt templates are used to evaluate the unlearning effectiveness, in-domain retainability, and cross-domain retainability of each method. Images with **green** frame denote desirable results, while the ones with **red** frame denote unlearning or retaining failures.





Figure A12: Visualization of the images generated by the unlearned DMs using different unlearning methods in the absence or presence of adversarial prompts.

## F Broader Use Cases of UNLEARNCANVAS

Although UNLEARNCANVAS is originally designed for benchmarking MU methods, we would like to demonstrate its broader use cases of benchmarking more generative modeling tasks, thanks to its good properties discussed in Sec. 3. In this section, we start with a case study on the task of style transfer, which is a much more well-studied topic than MU, but surprisingly also faces great challenges in building up a comprehensive and precise evaluation framework. In the next, we will first dissect the key challenges of the current style transfer evaluation framework, and then demonstrate how UNLEARNCANVAS efficiently resolves these challenges and further proposes a comprehensive and automated benchmark. Through extensive experiments, we draw demonstrate new insights from these results and illuminate the challenges of the future research directions. In the end, we will discuss the possibility of using UNLEARNCANVAS to benchmark more generative modeling tasks.

### F.1 Benchmarking Style Transfer using UNLEARNCANVAS

**Style transfer.** Style transfer is a long-standing topic and focuses on transferring the artistic style from one style image (also known as the *reference* image)  $\mathbf{x}_s$  to a target content image  $\mathbf{x}_c$ . The most recent methods typically employ a neural network, denoted by  $\theta_s$ , to extract style features and perform stylization in a single inference step, expressed as  $\hat{\mathbf{x}}_o = f_{\theta_s}(\mathbf{x}_s, \mathbf{x}_c)$ . Current state-of-the-art (SOTA) style transfer techniques exhibit remarkable generalization capabilities, successfully transferring styles not encountered during training and not requiring any further back-propagations. Figure 1 illustrates the pipeline of this task. Most existing literature [74–79] utilize WIKIART [66] to provide different styles for training the stylization network  $\theta_s$ . During the evaluation, the validation set of WIKIART will serve as the style (reference) images, together with the content images from the COCO dataset [80], to form a test bed for style transfer and style learning methods. However, such an evaluation scheme has some inherent limitations, which will be detailed below.

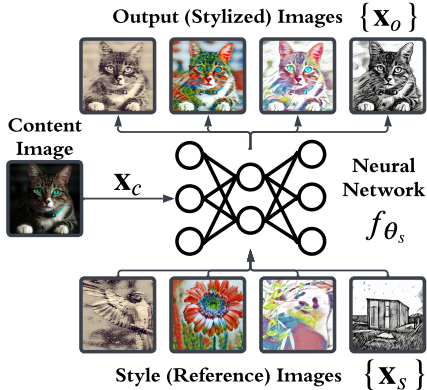


Figure A13: An illustration of the task of style transfer.

**Issues and challenges in the evaluating methods for style transfer** Although the task of style transfer has been widely studied, its evaluations are still based on very limited quantitative or even solely qualitative assessments, potentially leading to incomplete and inaccurate assessments [81]. Upon examining the evaluation pipelines of over 10 state-of-the-art (SOTA) style transfer methods, three significant challenges are identified within the current widely accepted evaluation frameworks.

- **Challenge I (C1): The lack of the ground truth images for style similarity evaluation.** Unlike other vision tasks like classification [82], detection [83], and segmentation [84], one of the key shortcomings of the current style transfer evaluation lies in the lack of the ground truth images  $\mathbf{x}_g$  for the given reference style image  $\mathbf{x}_s$  and the content image  $\mathbf{x}_c$ . Consequently, existing evaluation metrics, such as the style loss  $\ell_{\text{style}}$  [85, 86], has to be calculated with the reference style image  $\mathbf{x}_s$  as the ground truth  $\mathbf{x}_g$ , namely  $\ell_{\text{style}}(\mathbf{x}_s, \hat{\mathbf{x}}_o)$ , rather than directly using the ground truth  $\ell_{\text{style}}(\mathbf{x}_o, \mathbf{x}_g)$ . Obviously, such an indirect evaluation may lead to inaccurate results due to the different contents held in  $\mathbf{x}_s$  and  $\mathbf{x}_o$ . Existing work has demonstrated that such evaluation metrics can lead to very different conclusion from that of the user study [79]. Therefore, the creation of a *supervised* dataset with ground truth stylized images for every content image under each style is a timely remedy.
- **Challenge II (C2): The lack of algorithm stability evaluation against varied reference images  $\mathbf{x}_s$ .** The evaluation of algorithm stability in style transfer has often been neglected. This assessment requires consistent performance of a method on a content image  $\mathbf{x}_c$  across different style reference images  $\mathbf{x}_s$  representing the same target artistic style. Ideally, an algorithm should maintain uniform quality across various references within a style, avoiding significant performance variations. Current challenges in such evaluations stem from the lack of reference sets that exhibit *high stylistic consistency* within each style. Figure A1 showcases examples from the widely-used WIKIART dataset. Despite sharing the same artistic label, images within a row show considerable divergence in visual appearance and style. Consequently, using these images as style references can lead to stable algorithms producing stylistically varied outputs, leading to misleading assessment results. Therefore, creating a dataset

with high stylistic uniformity within each style category and clear differentiation between styles is crucial for accurate measurement of algorithm stability.

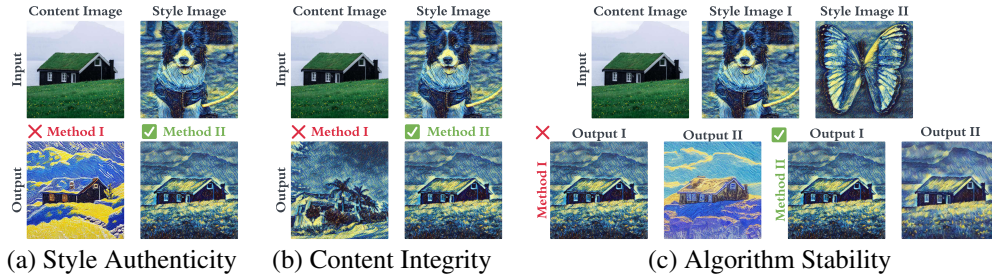


Figure A14: An illustration of comprehensive performance metrics for style transfer tasks. (a) Style authenticity measures the style similarity between the generated image and the style image. (b) Content integrity measures the content preservation between the generated image and the content image. (c) Algorithm stability reflects the sensitivity of the algorithm to different style images. For all the metrics in the illustration, “Method II” is always better than “Method I”.

### Building up a comprehensive evaluation pipeline with UNLEARNCANVAS for style transfer.

To revolutionize the evaluation framework with UNLEARNCANVAS, it is crucial to first understand the components that form a comprehensive evaluation pipeline, ensuring a thorough and impartial assessment of performance. In the realm of style transfer, three pivotal metrics stand paramount, as depicted in [Figure A14](#):

① **Style Authenticity**: Assesses the extent to which the style of the generated image aligns with that of the provided reference style image(s). ② **Content Integrity**: Measures how well the content features are preserved post style transfer. ③ **Algorithm Stability**: Evaluates the algorithm’s robustness against variations in content subjects, target styles, or style reference image selections, while keeping other parameters constant.

In addressing the challenges (C1-C2), UNLEARNCANVAS proves to be inherently advantageous. **First**, the style-specific supervision embedded in UNLEARNCANVAS enables the provision of ground truth for quantitative evaluations of stylized images. **Second**, the extensive image collection within the same style category in UNLEARNCANVAS allows for the assessment of algorithm stability by applying style transfer methods to varied reference images. To encompass the aforementioned aspects of style transfer performance, we propose the following quantitative evaluation metrics:

① **Style Loss**: Utilizes a feature map-based style loss [87] to quantify the stylistic *dissimilarity* between image pairs, effectively representing the inverse of style authenticity. ② **Content Loss**: Employs a VGG-based, feature-map content loss [87, 88] to measure the visual dissimilarity between the reference and generated images, essentially mirroring content integrity. ③ **Averaged Standard Deviation (STD)**: Computes the average STD of style and content loss *w.r.t.* the same reference image, reflecting algorithm stability.

Furthermore, similar to the MU task, we also consider efficiency metrics for each method, including:

④ **Average Time Consumption**: Measures the time required for performing style transfer. ⑤ **Peak GPU Memory Consumption**: Records the maximum GPU memory usage during the style transfer process. ⑥ **Model Storage Memory Consumption**: Assesses the memory requirement for storing the style transfer model.

With the introduction of evaluation metrics (①-⑥), we establish a comprehensive evaluation pipeline for style transfer. The process is as follows:

For each style transfer method, style transfer is executed within each object class. Specifically, for every style, images indexed from 1 to 18 serve as style reference images, while seed images corresponding to indices 19 and 20 are used as content images. Style and content loss are computed for each pair of style reference and content images, with the stylized images derived from the used seed content images serving as the ground truth for style loss. This results in a total of  $60 \times 20 \times 18 \times 2$  experimental trials for each method. For a specific style and content image pair, 18 experiments are performed to calculate the Standard Deviation (STD) values for both style and content loss. These

STD values are then averaged over all content images (amounting to  $60 \times 20 \times 2$  cases). The findings from these comprehensive evaluations are presented in Table A6.

Following this evaluation pipeline, we scrutinized 9 prominent style transfer methods, including SANET [78], MCC [89], MAST [90], ARTFLOW with its two variants (AF-ADAIn and AF-WCT) [79], IE-CONTRAST [91], CAST [88], STYTR2 [87], and BLIP [92].

Table A6: Performance overview of different style transfer methods evaluated with UNLEARNCANVAS dataset. The performance are assessed from the perspectives of stylistic authenticity (style loss), content integrity (content loss), algorithm stability (standard deviations from different dimensions), and efficiency. For all the metrics, *smaller* values are always preferred for better performance. The best performance per each metric is highlighted in **bold**. The standard deviations are first calculated with respect to different styles, object classes, or tested content images and then averaged in order to depict the algorithm stability from different perspectives.

Method	Style Loss STD (Averaged over)				Content Loss STD (Averaged over)				Efficiency		
	Mean	Style	Object	Content Image	Mean	Style	Object	Content Image	Time (s/image)	Memory (GB)	Storage (GB)
SANET	23.48	<b>2.73</b>	<b>2.87</b>	1.87	0.85	0.12	0.17	0.09	0.29	<b>2.3</b>	0.11
MCC	<b>17.92</b>	4.59	4.82	2.14	0.96	0.14	0.21	0.07	0.38	5.4	0.10
MAST	24.10	2.87	3.16	1.74	1.42	0.33	0.34	0.18	2.86	4.8	0.16
AF-ADAIn	20.78	2.96	3.13	1.65	1.09	0.11	0.19	0.05	0.53	6.3	0.08
AF-WCT	20.22	2.94	3.19	1.75	1.02	0.12	0.18	0.05	0.53	6.3	<b>0.08</b>
IE-CONTRAST	21.27	3.01	3.32	2.05	1.08	0.29	0.31	0.15	<b>0.05</b>	3.8	0.11
CAST	24.01	2.78	2.90	<b>1.35</b>	1.38	0.32	0.40	0.16	0.32	6.7	0.19
STYTR2	19.75	3.04	3.30	1.91	<b>0.62</b>	<b>0.10</b>	<b>0.12</b>	<b>0.04</b>	0.58	3.9	0.21
BLIP	25.43	2.90	3.06	2.03	1.61	0.30	0.34	0.16	8.87	7.2	7.23

**Experiment results analysis.** Tab. A6 provides a systematic evaluation of the performance of various methods tested. From the analysis, we can derive several crucial insights:

**First**, it is evident that no single method excels across all evaluation metrics. Notably, MCC demonstrates superior performance in maintaining stylistic authenticity, as indicated by a low style loss. Conversely, STYTR2 stands out in preserving content integrity, reflected by its minimal content loss.

**Second**, the assessment of standard deviation is indispensable for a comprehensive evaluation. The method with the optimal performance does not necessarily exhibit the greatest stability. This is particularly apparent in the style loss evaluation, where MCC, despite achieving the best result in terms of style loss, exhibits the least stability, denoted by the highest standard deviation.

## F.2 Other Possible Applications of UNLEARNCANVAS

In the preceding section, we demonstrated the application of UNLEARNCANVAS in refining evaluation metrics and frameworks for style transfer. Beyond this, we recognize the potential of UNLEARNCANVAS in diverse domains. Here, we delve into two illustrative examples:

**Bias mitigation.** Bias mitigation in DMs, which are now gaining popularity, can also benefit from UNLEARNCANVAS. Its hierarchical and balanced architecture enables the deliberate introduction of artificial biases by selectively omitting data from specific groups. For instance, by predominantly excluding images from styles other than the ‘Van Gogh Style’ within the ‘Dogs’ class, DMs finetuned on this dataset will inherently exhibit a tendency to generate images of dogs in the Van Gogh style, particularly when the style is not explicitly specified in the prompt. This approach not only allows for the manipulation and quantification of biases but also paves the way for UNLEARNCANVAS to become a standardized benchmark for bias mitigation, similar to the role of MU for DMs.

**Vision in-context learning (V-ICL).** V-ICL [93–96] is another domain where UNLEARNCANVAS can be effectively applied. The field of V-ICL is in urgent need of robust, comprehensive methods for the fair assessment of existing models. In this context, the image pairs from UNLEARNCANVAS are ideally suited for evaluating various tasks such as style transfer, image inpainting, and image segmentation, offering a rich resource for nuanced and quantitative analyses.

## G Impact Statement

This work helps improve the assessment and further promotes the advancement of MU (machine unlearning) methods for DMs (diffusion models), which are known to be effective in relieving or mitigating the various negative societal influences brought by the prevalent usage of DMs, which include but are not limited to the following aspects.

- **Avoiding Copyright Issues.** There is an urgent need for the generative model providers to scrub the influence of certain data on an already-trained model. In January 2023, a notable lawsuit targeted two leading AI art generators, Stable Diffusion [1] and Midjourney [2], for alleged copyright infringement. Concurrently, incidents with the recently released Midjourney V6 [2] also highlighted a visual plagiarism issue on famous film scenes. These instances illuminate the broad copyright challenges inherent in the way of training data collection method of those foundation generative models' training datasets. MU methods can be used as an effective method to remove the influence of the private data and avoid unnecessary retraining.

- **Mitigating biases and stereotypes.** Generative AI systems are known to have tendencies towards bias, stereotypes, and reductionism, when it comes to gender, race and national identities [17]. For example, a recent study on the images generated with Midjourney revealed, that images associated with higher-paying job titles featured people with lighter skin tones, and that results for most professional roles were male-dominated [16]. MU is known to be effective in eliminating biases rooted in the training data. Moreover, UNLEARNCANVAS offers a flexible framework to benchmark MU techniques against bias removal, allowing for the creation and quantitative control of biases across different object classes for comprehensive bias removal studies.