

# ArtiFade: Learning to Generate High-quality Subject from Blemished Images

Shuya Yang\*, Shaozhe Hao\*, Yukang Cao<sup>†</sup>, Kwan-Yee K. Wong<sup>†</sup>

The University of Hong Kong

## Abstract

Subject-driven text-to-image generation has witnessed remarkable advancements in its ability to learn and capture characteristics of a subject using only a limited number of images. However, existing methods commonly rely on high-quality images for training and may struggle to generate reasonable images when the input images are blemished by artifacts. This is primarily attributed to the inadequate capability of current techniques in distinguishing subject-related features from disruptive artifacts. In this paper, we introduce ArtiFade to tackle this issue and successfully generate high-quality artifact-free images from blemished datasets. Specifically, ArtiFade exploits fine-tuning of a pre-trained text-to-image model, aiming to remove artifacts. The elimination of artifacts is achieved by utilizing a specialized dataset that encompasses both unblemished images and their corresponding blemished counterparts during fine-tuning. ArtiFade also ensures the preservation of the original generative capabilities inherent within the diffusion model, thereby enhancing the overall performance of subject-driven methods in generating high-quality and artifact-free images. We further devise evaluation benchmarks tailored for this task. Through extensive qualitative and quantitative experiments, we demonstrate the generalizability of ArtiFade in effective artifact removal under both in-distribution and out-of-distribution scenarios.

## 1 Introduction

With the rapid advancement of generative diffusion models (Rombach et al. 2022; Song, Meng, and Ermon 2021; Saharia et al. 2022; Zhang, Rao, and Agrawala 2023; Ho, Jain, and Abbeel 2020), subject-driven text-to-image generation (Gal et al. 2023; Ruiz et al. 2023; Kumari et al. 2023; Kawar et al. 2023; Chen et al. 2023), which aims to capture distinct characteristics of a subject by learning from a few images of the subject, has gained significant attention. This approach empowers individuals to seamlessly incorporate their preferred subjects into diverse and visually captivating scenes by simply providing text conditions. Representative works such as Textual Inversion (Gal et al. 2023) and DreamBooth (Ruiz et al. 2023) have shown promising results on this task. Specifically, Textual Inversion proposes to optimize a textual embedding to encode identity characteristics that provide rich subject information for subsequent

generation. DreamBooth shares a similar idea but additionally fine-tunes the diffusion model to preserve more identity semantics. Plenty of successive efforts have been made to advance this task from various perspectives, including generation quality, compositionality, and efficiency (Kumari et al. 2023; Chen et al. 2023; Kawar et al. 2023).

Both of the above mentioned methods, along with their follow-up works, however, rely heavily on the presence of unblemished input images that contain only relevant identity information. This is often expensive or even unavailable in real-world applications. Instead, in practical scenarios such as scraping web images of a desired subject, it is common to encounter images that are blemished by various *visible* artifacts such as watermarks, drawings, and stickers. Additionally, there also exist *invisible* artifacts like adversarial noises (Van Le et al. 2023) that are not easily detectable or removable using off-the-shelf tools. These artifacts can significantly impede the comprehensive learning of the subject and lead to a catastrophic decline in performance across multiple dimensions (see Fig. 1). This limitation arises from the feature confusion inherent in the existing subject-driven learning process. The process simultaneously captures subject-related features and disruptive artifact interference. It lacks the discriminative power to distinguish these two from each other, and fails to preserve the integrity of subject characteristics while mitigating negative effects caused by artifacts. As blemished inputs are inevitable in applications, a pressing challenge emerges: **Can we effectively perform subject-driven text-to-image generation using blemished images?** We term this novel problem (*i.e.*, generating subject-driven images from blemished inputs) as blemished subject-driven generation in this paper.

To answer the above question, we present **ArtiFade**, the first model to tackle blemished subject-driven generation by adapting vanilla subject-driven methods (*e.g.*, Textual Inversion (Gal et al. 2023) and DreamBooth (Ruiz et al. 2023)) to effectively extract subject-specific information from blemished training data. The key objective of ArtiFade is to learn the implicit relationship between natural images and their blemished counterparts through alignment optimization. Specifically, we introduce a specialized dataset construction method to create pairs of unblemished images and their corresponding counterparts. These pairs can be applied to fine-tune various subject-driven approaches in the context

\*Equal contribution

<sup>†</sup>Corresponding authors

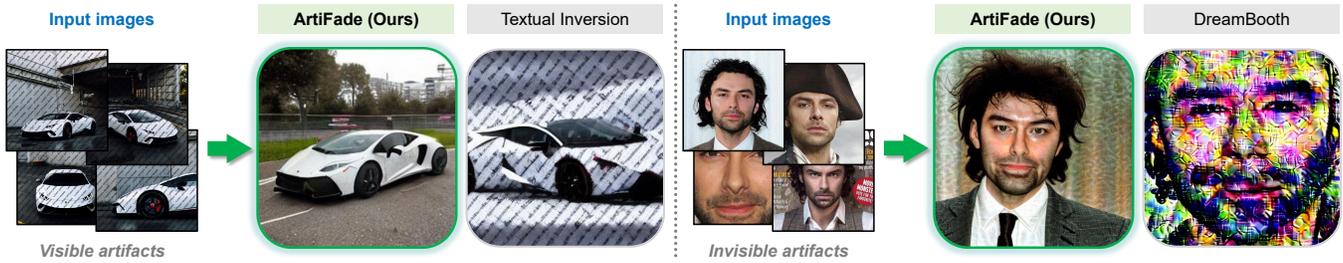


Figure 1: Blemished subject-driven generation with our ArtiFade and vanilla subject-driven methods. We display images generated using ArtiFade and Textual Inversion on watermark artifacts on the left, and ArtiFade and DreamBooth on adversarial noise artifacts (Van Le et al. 2023) on the right. In contrast to the poor performance of Textual Inversion and DreamBooth, which are negatively affected by the visible or invisible artifacts, ArtiFade produces much better fidelity of the subject with high-quality generation.

of blemished subject-driven generation. Besides, we also observe fine-tuning an extra learnable embedding in the textual space, named artifact-free embedding, can enhance prompt fidelity in the blemished subject-driven generation.

We further introduce an evaluation benchmark that encompasses (1) multiple test sets of blemished images with diverse artifacts, and (2) tailored metrics for accurately assessing the performance of blemished subject-driven generation methods. A thorough experimental evaluation shows that our method consistently outperforms other existing methods, both qualitatively and quantitatively. Notably, ArtiFade exhibits superb capabilities in handling out-of-distribution (OOD) scenarios involving diverse types of artifacts that are distinct from the training data. This inherent generalizability indicates our model can effectively learn to discern and distinguish the patterns exhibited by artifacts and unblemished images, instead of overfitting to a specific type of artifacts.

In summary, our key contributions are as follows:

- We are the first to tackle the novel challenge of blemished subject-driven generation. To address this task, we propose ArtiFade that fine-tunes diffusion models to align unblemished and blemished data.
- We introduce an evaluation benchmark tailored for effectively assessing the performance of blemished subject-driven generation techniques.
- We conduct extensive experiments and demonstrate that ArtiFade outperforms current methods significantly. We show noteworthy generalizability of ArtiFade, effectively addressing both in-distribution and out-of-distribution scenarios with various types of artifacts.

## 2 Related Work

**Text-to-image synthesis** Text-to-image generation has attracted considerable attention in recent years by leveraging Generative Adversarial Networks (GANs) (Goodfellow et al. 2014) and diffusion models (Ho, Jain, and Abbeel 2020; Rombach et al. 2022). Reed *et al.* (2016) was the first to integrate GANs into text-to-image generation. Since then, several influential works had been proposed (Zhang et al. 2017, 2018; Xu et al. 2018; Zhang, Xie, and Yang

2018; Zhu et al. 2019; Li et al. 2019a; Ruan et al. 2021; Zhang et al. 2021; Cheng et al. 2020; Qiao et al. 2019; Yin et al. 2019), demonstrating impressive results with improved resolution (Zhang et al. 2017, 2018) and fidelity of fine details (Xu et al. 2018). Diffusion models in text-to-image synthesis have also yielded remarkable results owing to their ability in generating precise and customized images that better align with individual text specifications (Nichol et al. 2022; Saharia et al. 2022; Ramesh et al. 2022; Gu et al. 2022; Rombach et al. 2022).

**Subject-driven generation** Subject-driven generation has gained popularity due to its ability to generate personalized images based on a given set of subject images and text prompts. One prominent method in subject-driven generation is Textual Inversion (Gal et al. 2023), which involves learning an embedding vector by minimizing the Latent Diffusion Model loss (Rombach et al. 2022) on input images. The learned embedding vector can be effectively combined with text prompts, allowing seamless integration in the text-to-image generation process. Recent approaches (Ruiz et al. 2023; Kumari et al. 2023; Lu et al. 2023) have significantly enhanced subject reconstruction fidelity by incorporating fine-tuning techniques.

**Artifacts removal** Shadow and watermark removal are classic tasks in image processing and computer vision. At the early stage, most approaches for shadow removal or image recovery relied on the properties of intensity and illumination (Finlayson, Drew, and Lu 2009; Finlayson et al. 2006; Zhang, Zhang, and Xiao 2015; Xiao et al. 2013b,a; Finlayson, Hordley, and Drew 2002; Khan et al. 2015; Shor and Lischinski 2008; Arbel and Hel-Or 2010; Guo, Dai, and Hoiem 2011). Some methods also incorporated color features to improve their results (Guo, Dai, and Hoiem 2011). Deep learning techniques and Convolutional Neural Networks (CNNs) have played a significant role in advancing shadow removal methods and producing impressive results in recent years (Ding et al. 2019; Hu et al. 2019; Le and Samaras 2019; Liu et al. 2021; Wang, Li, and Yang 2018; Zhu et al. 2022; Chen et al. 2021; Jin et al. 2023; Fu et al. 2021). Several studies (Wang, Li, and Yang 2018; Liu et al. 2021; Hu et al. 2019; Ding et al. 2019) have incor-

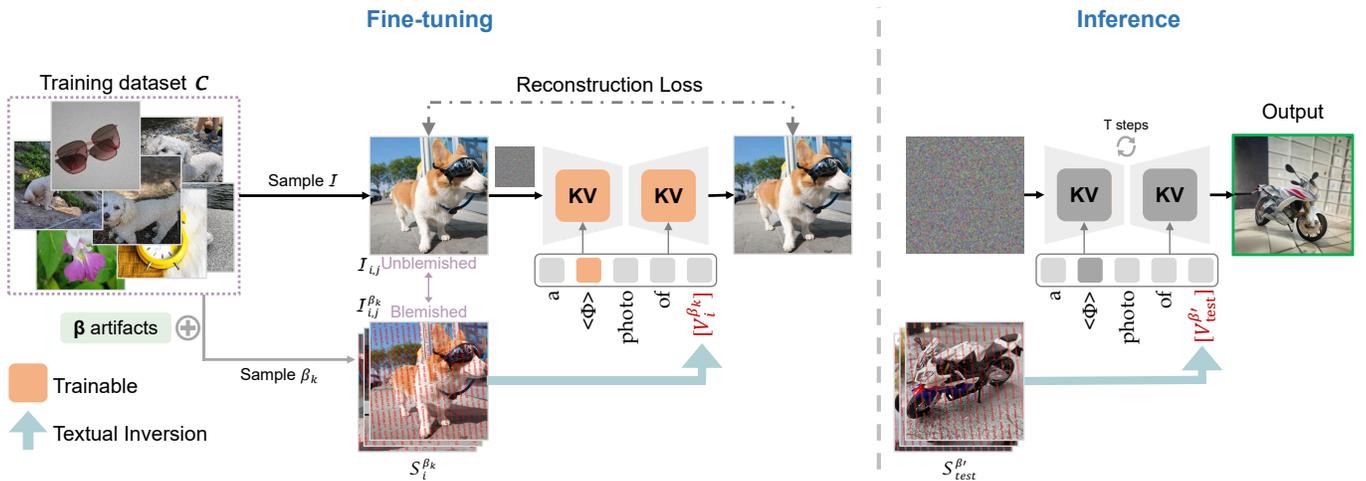


Figure 2: Overview of ArtiFade. On the left, we present Artifact Rectification Training, which involves an iterative process of calculating reconstruction loss between an unblemished image and the reconstruction of its blemished embedding. The right-hand side is the inference stage that tests ArtiFade on unseen blemished images. To avoid ambiguity, we (1) simplify the training of Textual Inversion into an input-output form, and (2) use “fine-tuning” and “inference” to respectively refer to the fine-tuning stage of ArtiFade and the use of ArtiFade for subject-driven generation.

porated GANs to further enhance the results of shadow removal techniques. Moreover, with the increasing popularity of diffusion models in image generation, a novel diffusion-based method for shadow removal has recently been introduced (Guo et al. 2023).

The most widely adopted methods for recovering concealed information from watermarked images include the application of generalized multi-image matting algorithms (Dekel et al. 2017), complemented by image inpainting techniques (Xu, Lu, and Zhou 2017; Qin et al. 2018; Huang and Wu 2004), and the utilization of deep neural networks and CNNs (Cheng et al. 2018). Similar to shadow removal, GANs and Conditional GANs (Mirza and Osindero 2014) are also widely used in watermark removal tasks (Li et al. 2019b; Cao et al. 2019; Liu, Zhu, and Bai 2021). Our work is closely related to these previously mentioned studies. We are the first to address the artifact issues in the realm of subject-driven text-to-image generation.

### 3 Method

Given a set of blemished input images, our objective is to eliminate their negative impacts on the quality of subject-driven image generation. To achieve this goal, we present ArtiFade, an efficient framework that learns to discern and distinguish the patterns exhibited by various types of artifacts and unblemished images. In this section, we focus exclusively on ArtiFade based on Textual Inversion. However, it is important to note that the ArtiFade framework can be generalized to other subject-driven generation methods. As shown in Fig. 2, ArtiFade based on Textual Inversion incorporates two main components, namely the fine-tuning of the partial parameters (*i.e.*, key and value weights) in the diffusion model and the simultaneous optimization of an artifact-free embedding  $\langle \Phi \rangle$ . We begin by discussing the preliminaries of the Latent Diffusion Model and Textual Inversion.

In the following subsections, we elaborate our automatic construction of the training dataset, which consists of both blemished and unblemished data, illustrated in Sec. 3.1. We then introduce Artifact Rectification Training, a method for fine-tuning the model to accommodate blemished images, as discussed in Sec. 3.2. We finally present the use of ArtiFade for handling blemished images in Sec. 3.3.

**Preliminary** Latent Diffusion Model (LDM) (Rombach et al. 2022) is a latent text-to-image diffusion model derived from Diffusion Denoising Probabilistic Model (DDPM) (Ho, Jain, and Abbeel 2020). LDM leverages a pre-trained autoencoder to map image features between the image and latent space. This autoencoder comprises an encoder  $\mathcal{E}$ , which transforms images into latent representations, and a decoder  $\mathcal{D}$ , which converts latent representations back into images. The autoencoder is optimized using a set of images so that the reconstructed image  $\hat{x} \approx \mathcal{D}(\mathcal{E}(x))$ . Additionally, LDM introduces cross-attention layers (Vaswani et al. 2017) within the U-Net (Ronneberger, Fischer, and Brox 2015), enabling the integration of text prompts as conditional information during the image generation process. The LDM loss is defined as

$$\mathcal{L}_{LDM} := \mathbb{E}_{z \sim \mathcal{E}(\mathcal{I}), y, \epsilon \sim N(0,1)} \left[ \|\epsilon - \epsilon_{\theta}(z_t, t, y)\|_2^2 \right], \quad (1)$$

where  $\mathcal{E}$  encodes the image  $\mathcal{I}$  into the latent representation  $z$ . Here,  $z_t$  denotes the noisy latent representation at timestep  $t$ ,  $\epsilon_{\theta}$  refers to the denoising network, and  $y$  represents the text condition that is passed to the cross-attention layer.

Based on LDM, Textual Inversion (Gal et al. 2023) aims to capture the characteristics of a specific subject from a small set of images. Specifically, Textual Inversion learns a unique textual embedding by minimizing Eq. (1) on a few images that contain the particular subject. It can produce promising generation results with high-quality inputs, but

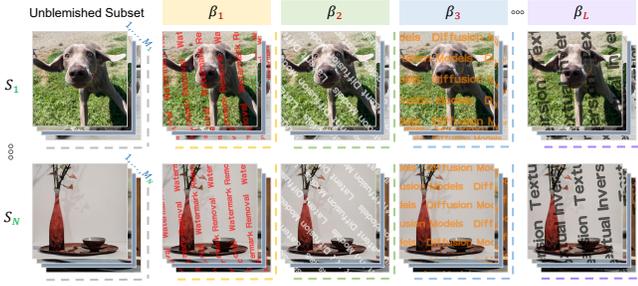


Figure 3: Examples of training dataset  $\mathcal{D}$  that contains both unblemished images and blemished counterparts.

fails on input images that are blemished by artifacts (see Fig. 1). This problem arises from the inherent limitation of Textual Inversion in learning shared characteristics exhibited in the input images without the capability in differentiating artifacts from unblemished subjects. In this paper, we aim to address this issue on deteriorated generation quality of Textual Inversion in the presence of blemished images.

### 3.1 Dataset Construction

Existing subject-driven generation methods operate under the assumption of unblemished training data, consisting of solely high-quality images devoid of any artifacts. However, this assumption does not align with real-world applications, where obtaining blemished images from the internet is a commonplace. To address this blemished subject-driven generation in this paper, we first construct a training set that incorporates both unblemished images and their blemished counterparts that are augmented with artifacts.

**Augmentation of multiple artifacts** We construct our dataset by collecting a multi-subject set  $\mathcal{C}$  of  $N$  image subsets from existing works (Gal et al. 2023; Ruiz et al. 2023; Kumari et al. 2023) and a set  $\mathcal{B}$  of  $L$  different artifacts:

$$\mathcal{C} = \{\mathcal{S}_i\}_{i=1}^N, \quad \mathcal{S}_i = \{\mathcal{I}_{i,j}\}_{j=1}^{M_i}, \quad \mathcal{B} = \{\beta_k\}_{k=1}^L, \quad (2)$$

where  $\mathcal{S}_i$  denotes the image subset corresponding to the  $i$ th subject,  $M_i$  is the total number of images in  $\mathcal{S}_i$ , and  $\beta_k$  represents a type of artifact for image augmentation. Our dataset  $\mathcal{D}$  can then be constructed by applying each artifact  $\beta_k$  to each image  $\mathcal{I}$  in  $\mathcal{S}_i$  separately, i.e.,

$$\mathcal{S}_i^{\beta_k} = \{\mathcal{I}_{i,j}^{\beta_k}\}_{j=1}^{M_i}, \quad \mathcal{D} = \{\mathcal{S}_i, \{\mathcal{S}_i^{\beta_k}\}_{k=1}^L\}_{i=1}^N, \quad (3)$$

where  $\mathcal{I}_{i,j}^{\beta_k}$  is the counterpart of  $\mathcal{I}_{i,j}$  augmented with the specific artifact  $\beta_k$ . Some examples of original images and their augmented versions with distinct artifacts can be found in Fig. 3. See the Appendix for more visualizations.

**Blemished textual embedding** For each blemished subset, we perform Textual Inversion to optimize a blemished textual embedding  $[v_i^{\beta_k}]$ , i.e.,

$$\mathcal{S}_i^{\beta_k} \xrightarrow{\text{Textual Inversion}} [v_i^{\beta_k}], \quad (4)$$

$i = 1, 2, \dots, N; \quad k = 1, 2, \dots, L$

By applying Eq. (4) on  $N$  subsets with  $L$  types of artifacts, we end up with a set of  $N \times L$  blemished textual embeddings  $\mathcal{V} = \{[v_i^{\beta_k}]\}_{i=1, k=1}^{N, L}$ , which will be used in the subsequent model fine-tuning. As we have illustrated in Fig. 1, directly prompting the diffusion model with  $[v_i^{\beta_k}]$  will lead to a significant decrease in generation quality. Consequently, our objective is to robustly handle blemished embeddings and effectively eliminate the detrimental impact of artifacts. We achieve this by devising a partial fine-tuning paradigm for the pre-trained diffusion model on the constructed training set  $\mathcal{D}$ , as elaborated in the following subsection.

### 3.2 Artifact Rectification Training

After establishing the curated dataset  $\mathcal{D}$ , we embark on training a generalizable framework on  $\mathcal{D}$ , capable of generating unblemished images using blemished textual embeddings. To this end, we propose artifact rectification training, which consists of two key components, namely partial fine-tuning of a pre-trained diffusion model and the optimization of an artifact-free embedding, to eliminate the artifacts and distortions in the generated images.

We fine-tune only partial parameters that are involved in processing the textual conditions. This strategy allows us to optimize the relevant components associated with the blemished textual embedding  $[v_i^{\beta_k}]$ . Considering that only the key and value weights in the diffusion model’s cross-attention layer are involved in the processing of textual embedding, we choose to fine-tune these two types of parameters  $W^k$  and  $W^v$ . Moreover, we find that optimizing an additional embedding,  $\langle \Phi \rangle$ , in the textual space with partial parameters could improve prompt fidelity by retaining the textual information of the model, as presented later in Sec. 4.7.

**Training objective** During each iteration, we will first randomly sample an unblemished image  $\mathcal{I}_{i,j}$  from the training set  $\mathcal{D}$  and a type of artifact  $\beta_k \in \mathcal{B}$  to obtain the blemished textual embedding  $[v_i^{\beta_k}] \in \mathcal{V}$  that is optimized on the blemished subset  $\mathcal{S}_i^{\beta_k}$ .

Specifically, given the sampled blemished textual embedding  $[v_i^{\beta_k}]$ , we form the prompt “a  $\langle \Phi \rangle$  photo of  $[v_i^{\beta_k}]$ ”, which will be input to the text encoder to acquire the text condition  $y_i^{\beta_k}$ . Our optimization objective will then be defined as reconstructing the unblemished image  $\mathcal{I}_{i,j}$  by conditioning the denoising process on the text condition  $y_i^{\beta_k}$ . Thus, we can formulate the final loss for training ArtiFade as

$$\mathcal{L}_{\text{ArtiFade}} := \mathbb{E}_{z \sim \mathcal{E}(\mathcal{I}_{i,j}), y_i^{\beta_k}, \epsilon \sim \mathcal{N}(0,1)} \left[ \|\epsilon - \epsilon_{\{W^k, W^v, \langle \Phi \rangle\}}(z_t, t, y_i^{\beta_k})\|_2^2 \right], \quad (5)$$

where  $\{W^k, W^v, \langle \Phi \rangle\}$  is the set of the trainable parameters of ArtiFade.

### 3.3 Subject-driven Generation with Blemished Images

After artifact rectification training, we obtain the ArtiFade model, prepared for the task of blemished subject-driven

Method	$\overline{\text{WM-model}}$ on WM-ID-test				
	I <sup>DINO</sup> ↑	R <sup>DINO</sup> ↑	I <sup>CLIP</sup> ↑	R <sup>CLIP</sup> ↑	T <sup>CLIP</sup> ↑
TI (unblemished)	0.488	1.349	0.730	1.070	0.283
TI (blemished)	0.217	0.852	0.576	0.909	0.263
Ours (TI-based)	<b>0.337</b>	<b>1.300</b>	<b>0.649</b>	<b>1.020</b>	<b>0.282</b>

Table 1: Quantitative results - ID.

Method	$\overline{\text{WM-model}}$ on WM-OD-test				
	I <sup>DINO</sup> ↑	R <sup>DINO</sup> ↑	I <sup>CLIP</sup> ↑	R <sup>CLIP</sup> ↑	T <sup>CLIP</sup> ↑
TI (unblemished)	0.488	1.278	0.730	1.136	0.283
TI (blemished)	0.229	0.858	0.575	0.929	0.262
Ours (TI-based)	<b>0.356</b>	<b>1.237</b>	<b>0.654</b>	<b>1.079</b>	<b>0.282</b>

Table 2: Quantitative results - OOD.

generation. Given a test image set  $S_{\text{test}}^{\beta'}$  in which all images are blemished by an arbitrary artifact  $\beta'$ , the ArtiFade model can generate high-quality subject-driven images using blemished samples with ease.

Specifically, we first obtain the blemished textual embedding  $[V_{\text{test}}^{\beta'}]$  by applying Textual Inversion on the test set  $S_{\text{test}}^{\beta'}$ . We then simply infer the ArtiFade model with a given text prompt that includes the blemished textual embedding, *i.e.*, “a  $\langle \Phi \rangle$  photo of  $[V_{\text{test}}^{\beta'}]$ ”. At the operational level, the sole distinction between our approach and vanilla Textual Inversion lies in inputting text prompts containing  $[V_{\text{test}}^{\beta'}]$  into the fine-tuned ArtiFade instead of the pre-trained diffusion model. This simple yet effective method resolves the issue of Textual Inversion’s incapacity to handle blemished input images, bearing practical utility.

**Details of ArtiFade models** We choose  $N = 20$  subjects, including pets, plants, containers, toys, and wearable items to ensure a diverse range of categories. We experiment with the ArtiFade model based on Textual Inversion trained with visible watermark artifacts, namely  $\overline{\text{WM-model}}$ . The training set of  $\overline{\text{WM-model}}$  involves  $L_{\text{WM}} = 10$  types of watermarks, characterized by various fonts, orientations, colors, sizes, and text contents. Therefore, we obtain 200 blemished subsets in total within the training set of  $\overline{\text{WM-model}}$ . We fine-tune  $\overline{\text{WM-model}}$  for a total of 16k steps.

## 4 Experiment

### 4.1 Implementation Details

We employ the pre-trained LDM (Rombach et al. 2022) following the official implementation of Textual Inversion (Gal et al. 2023) as our base diffusion model. We train the blemished textual embeddings for 5k steps using Textual Inversion. We use a learning rate of  $5e-3$  to optimize our Artifact-free embedding and  $3e-5$  for the partial fine-tuning of key and value weights. Note that all other parameters within the pre-trained diffusion model remain frozen. All experiments are conducted on 2 NVIDIA RTX 3090 GPUs. In the main paper, we focus on the comparison with Textual Inversion



Figure 4: Qualitative Comparison - ID. Unlike Textual Inversion which struggles to produce reasonable generation from blemished inputs, our method ( $\overline{\text{WM-model}}$ ) consistently learns the distinguished features of the given subject and achieves high-quality generation without distortion.

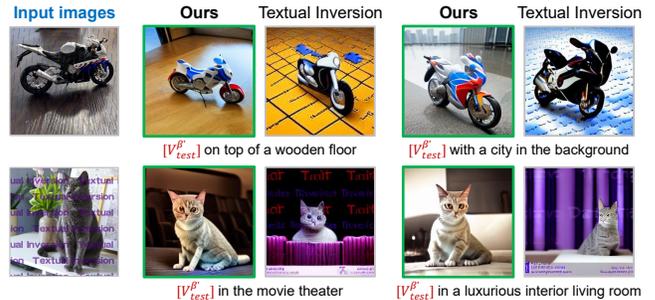


Figure 5: Qualitative Comparison - OOD. Our method ( $\overline{\text{WM-model}}$ ) is generalizable to process out-of-distribution artifacts that are unseen during the fine-tuning, demonstrating much better performance than Textual Inversion. Best viewed in PDF with zoom.

and DreamBooth to demonstrate the efficiency of our proposed contributions. See the Appendix for additional comparisons and applications.

### 4.2 Evaluation Benchmark

**Test dataset** We construct the test dataset using 16 novel subjects that differ from the subjects in the training set. These subjects encompass a wide range of categories, including pets, plants, toys, transportation, furniture, and wearable items. We form the visible test artifacts into two categories: (1) in-distribution watermarks (WM-ID-test) containing the same type as the training data, and (2) out-of-distribution watermarks (WM-OD-test) of different types from the training data. Within the WM-ID-test and WM-OD-test, we synthesize 5 distinct artifacts for each category, resulting in 80 test sets.

**Evaluation metrics** We evaluate the performance of blemished subject-driven generation from three perspectives: (1) the fidelity of subject reconstruction, (2) the fidelity of text conditioning, and (3) the effectiveness of mitigating the negative impacts of artifacts. Following common practice (Gal et al. 2023; Ruiz et al. 2023), we use CLIP (Radford et al. 2021) and DINO (Caron et al. 2021) similarities

for measuring these metrics. For the first metric, we calculate the CLIP and DINO similarity between the generated images and the unblemished version of the input images, respectively denoted as  $I^{\text{CLIP}}$  and  $I^{\text{DINO}}$ . For the second metric, we calculate the CLIP similarity between the generated images and the text prompt, denoted as  $T^{\text{CLIP}}$ . For the third metric, we calculate the relative ratio of similarities between generated images and unblemished input images compared to their blemished versions, defined as

$$R^{\text{CLIP}} = I^{\text{CLIP}}/I_{\beta}^{\text{CLIP}} \quad R^{\text{DINO}} = I^{\text{DINO}}/I_{\beta}^{\text{DINO}} \quad (6)$$

where  $I_{\beta}^{\text{CLIP}}$  and  $I_{\beta}^{\text{DINO}}$  respectively denote CLIP and DINO similarities between the generated images and the *blemished* input images. A relative ratio greater than 1 indicates that generated images resemble unblemished images more than blemished counterparts, suggesting fewer artifacts. Conversely, a ratio less than 1 indicates that generated images are heavily distorted with more artifacts. We use DINO ViT-S/16 (Caron et al. 2021) and CLIP ViT-B/32 (Radford et al. 2021) to compute all metrics.

### 4.3 Quantitative Comparisons

We conduct both in-distribution and out-of-distribution quantitative evaluations of our method and compare it to Textual Inversion with blemished embeddings. We additionally report the results using Textual Inversion on unblemished images as a reference, although it is not a direct comparison to our model.

**In-distribution (ID) analysis** We consider the in-distribution scenarios by testing  $\mathbb{W}\mathbb{M}\text{-model}$  on  $\mathbb{W}\mathbb{M}\text{-ID-test}$ . In Tab. 1, we can observe that the use of blemished embeddings in Textual Inversion leads to comprehensive performance decline including: (1) lower subject reconstruction fidelity (*i.e.*,  $I^{\text{DINO}}$  and  $I^{\text{CLIP}}$ ) due to the subject distortion in image generation; (2) lower efficiency for artifact removal (*i.e.*,  $R^{\text{DINO}}$  and  $R^{\text{CLIP}}$ ) due to inability to remove artifacts; (3) lower prompt fidelity (*i.e.*,  $T^{\text{CLIP}}$ ) since the prompt-guided background is unrecognizable due to blemishing artifacts. In contrast, our method consistently achieves higher scores than Textual Inversion with blemished embeddings across the board, demonstrating the efficiency of ArtiFade in various aspects.

**Out-of-distribution (OOD) analysis** We pleasantly discover that  $\mathbb{W}\mathbb{M}\text{-model}$  possesses the capability to handle out-of-distribution scenarios, owing to its training with watermarks of diverse types. We consider the out-of-distribution (OOD) scenarios for  $\mathbb{W}\mathbb{M}\text{-model}$  by testing it on  $\mathbb{W}\mathbb{M}\text{-OOD-test}$ , as presented in Tab. 2. Similar to ID evaluation, all of our metrics yield higher results than Textual Inversion with blemished embeddings. These results further demonstrate the generalizability of our method.

### 4.4 Qualitative Comparisons

We present qualitative comparisons between the output generated via ArtiFade and Textual Inversion with blemished textual embeddings, including in-distribution scenarios in Fig. 4 and out-of-distribution scenarios in Fig. 5.

Method	WM-ID-test				
	$I^{\text{DINO}}\uparrow$	$R^{\text{DINO}}\uparrow$	$I^{\text{CLIP}}\uparrow$	$R^{\text{CLIP}}\uparrow$	$T^{\text{CLIP}}\uparrow$
TI (unblemished)	0.488	1.349	0.730	1.070	0.283
TI (blemished)	0.217	0.852	0.576	0.909	0.263
DB (blemished)	0.503	0.874	0.738	0.939	0.272
Ours (TI-based)	0.337	1.300	0.649	1.020	0.282
Ours (DB-based)	<b>0.589</b>	<b>1.308</b>	<b>0.795</b>	<b>1.083</b>	<b>0.284</b>

Table 3: Quantitative comparison with DreamBooth.



Figure 6: Qualitative comparison with DreamBooth.

**In-distribution analysis** The images generated by Textual Inversion exhibit noticeable limitations when using blemished textual embeddings. Specifically, as depicted in Fig. 4, all rows predominantly exhibit cases of incorrect backgrounds that are highly polluted by watermarks. By using ArtiFade, we are able to eliminate the background watermarks.

**Out-of-distribution analysis** In addition, we conduct experiments with our  $\mathbb{W}\mathbb{M}\text{-model}$  to showcase its capability to remove out-of-distribution watermarks, as shown in Fig. 5. It is important to note that in the first row, the watermark in the input images may not be easily noticed by human eyes upon initial inspection due to the small font size and high image resolution. However, these artifacts have a significant effect when used to train blemished embeddings for generating images. ArtiFade effectively eliminates the artifacts on the generated images, improving reconstruction fidelity and background accuracy, hence leading to substantial enhancements in overall visual quality.

### 4.5 ArtiFade with DreamBooth

The ArtiFade fine-tuning framework is not limited to Textual Inversion with textual embedding; it can also be generalized to DreamBooth. We use the same training dataset and blemished subsets as in the case of the  $\mathbb{W}\mathbb{M}\text{-model}$  (*i.e.*,  $N = 20$ ,  $L_{WM} = 10$ ). The vanilla DreamBooth fine-tunes the whole UNet model, which conflicts with the fine-tuning parameters of ArtiFade. We therefore use DreamBooth with low-rank approximation (LoRA)<sup>1</sup> to train LoRA adapters (Hu et al. 2022) for the text encoder, value, and query weights of the diffusion model for each blemished subset using Stable Diffusion v1-5. For simplicity, we will

<sup>1</sup>[https://huggingface.co/docs/peft/main/en/task\\_guides/dreambooth\\_lora](https://huggingface.co/docs/peft/main/en/task_guides/dreambooth_lora)



Figure 7: Qualitative Comparison between ours and DreamBooth when inputs are blemished by invisible adversarial noises.

use DreamBooth to refer to DreamBooth with LoRA below. During the fine-tuning of DreamBooth-based ArtiFade, we load the pre-trained adapters and only unfreeze key weights since value weights are reserved for DreamBooth subject information. In Tab. 3, it is evident that our method, based on DreamBooth, yields the highest scores among all cases. Our method also maintains DreamBooth’s advantages in generating images with higher subject fidelity and more accurate text prompting, outperforming ArtiFade with Textual Inversion. We show some qualitative results in Fig. 6.

#### 4.6 Invisible Artifacts Blemished Subject Generation

ArtiFade demonstrates exceptional performance in handling subjects characterized by intricate features and blemished by imperceptible artifacts. We collect 20 human figure datasets from the VGGFace2 dataset (Cao et al. 2018). We then use the Anti-DreamBooth (Van Le et al. 2023) ASPL method to add adversarial noises to each group of images, producing 20 blemished datasets for fine-tuning a DreamBooth-based ArtiFade model. The model is fine-tuned for 12k steps. As illustrated in Fig. 7, our approach surpasses the DreamBooth in differentiating the learning of adversarial noises from human face features. In contrast to DreamBooth, which is fooled into overfitting adversarial noises, thereby generating images with a heavily polluted background, our model reconstructs human figures in image generation while maintaining high fidelity through text prompting.

#### 4.7 Ablation Studies

We conduct ablation studies to demonstrate the efficiency of our method by comparing with three alternative variants, which encompass (1)  $\text{Var}_A$ , where we solely fine-tune the artifact-free embedding; (2)  $\text{Var}_B$ , where we fine-tune parameters related to image features, *i.e.*, query weights  $W^q$ , along with the artifact-free embedding, and (3)  $\text{Var}_C$ , where we fine-tune key and value weights, *i.e.*,  $W^k$  and  $W^v$ , exclusively. We use our `WM-model` to compare it with other variants by testing on the `WM-ID-test`.

**Effect of partial fine-tuning** As shown in Tab. 4, compared to  $\text{Var}_A$ , our full method yields higher scores on all metrics by a significant margin, except for  $R^{\text{DINO}}$ . This

Method	$W^{kv}$	$W^q$	$\langle \Phi \rangle$	$I^{\text{DINO}}$	$R^{\text{DINO}}$	$I^{\text{CLIP}}$	$R^{\text{CLIP}}$	$T^{\text{CLIP}}$
$\text{Var}_A$			✓	0.154	<b>1.412</b>	0.566	0.984	0.265
$\text{Var}_B$		✓	✓	0.283	1.230	0.617	0.978	0.277
$\text{Var}_C$	✓			<b>0.342</b>	1.292	<b>0.652</b>	1.019	0.280
Ours	✓		✓	0.337	1.300	0.649	<b>1.020</b>	<b>0.282</b>

Table 4: Quantitative comparison of ablation studies.

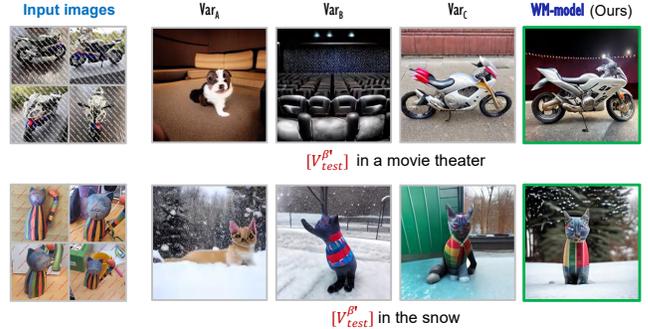
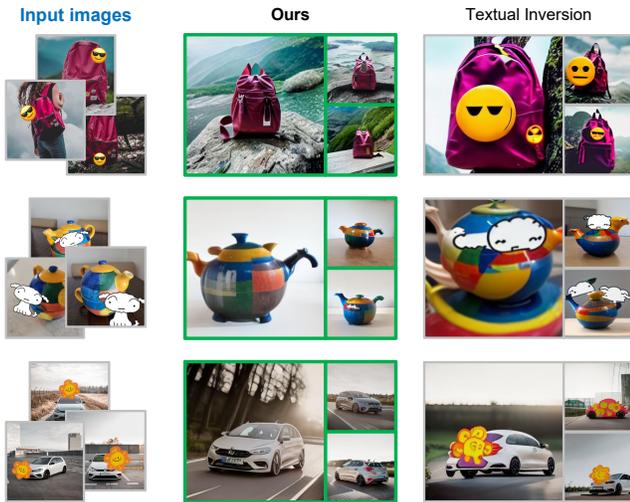


Figure 8: Qualitative comparison of ablation studies.

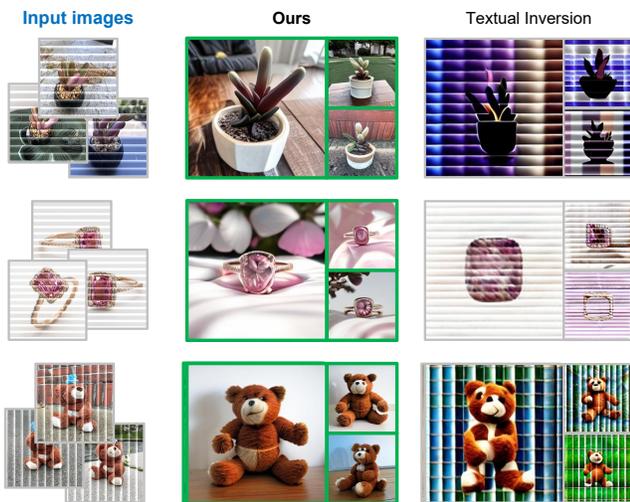
is reasonable, as the artifact-free embedding can be easily overfitted to the training data, resulting in generated images that resemble a fusion of training images (Fig. 8,  $\text{Var}_A$ ). As a result, the denominator of  $R^{\text{DINO}}$ , namely the similarity between the generated image and the blemished image, is significantly decreased, leading to a high  $R^{\text{DINO}}$ . Due to similar reason,  $\text{Var}_A$  shows lowest  $I^{\text{DINO}}$ ,  $I^{\text{CLIP}}$ , and  $T^{\text{CLIP}}$  among all variants, indicating that it fails to reconstruct the correct subject. Overall, both quantitative and qualitative evaluation showcases that solely optimizing the artifact-free embedding is insufficient to capture the distinct characteristics presented in the blemished input image, demonstrating the necessity of partial fine-tuning.

**Effect of fine-tuning key and value weights** As shown in Tab. 4 and Fig. 8,  $\text{Var}_B$  yields unsatisfactory outcomes in all aspects compared to ours. The lower  $R^{\text{DINO}}$  and  $R^{\text{CLIP}}$  suggest that the generated images retain artifact-like features and bear closer resemblances to the blemished subsets. Furthermore, the reduced  $T^{\text{CLIP}}$  indicates diminished prompt fidelity, as the approach fails to accurately reconstruct the subject from the blemished embeddings, which is also evidenced by Fig. 8. These findings suggest that fine-tuning the parameters associated with text features yields superior enhancements in terms of artifact removal and prompt fidelity.

**Effect of the artifact-free embedding** With  $\text{Var}_C$ , we exclude the optimization of artifact-free embedding. In Tab. 4, we can observe that  $\text{Var}_C$  yields higher  $I^{\text{DINO}}$  and  $I^{\text{CLIP}}$  but lower  $R^{\text{DINO}}$  and  $R^{\text{CLIP}}$  compared to our `WM-model`, which indicates that the approach achieves higher subject fidelity but lower efficiency in eliminating artifacts when generating images. Since our primary objective is to generate artifact-free images from blemished textual embedding, our `WM-model` chooses to trade off subject reconstruction fidelity for the ability to remove artifacts. Additionally, this



(a) Sticker removal.



(b) Glass effect removal.

Figure 9: Applications. Our  $\mathbb{W}\mathbb{M}\text{-model}$  can be applied to remove various unwanted artifacts in the input images, *e.g.* stickers, glass effect, etc.

approach produces lower  $T^{\text{CLIP}}$  than ours, suggesting that the artifact-free embedding effectively improves the model’s capability to better preserve text information (see Fig. 8).

## 5 More Applications

We apply our  $\mathbb{W}\mathbb{M}\text{-model}$  to more artifact cases, such as stickers and glass effects, showcasing its broad applicability.

**Sticker removal.** In Fig. 9a, we test  $\mathbb{W}\mathbb{M}\text{-model}$  on input images that are blemished by cartoon stickers. The cartoon sticker exhibits randomized dimensions and is positioned arbitrarily within each image.  $\mathbb{W}\mathbb{M}\text{-model}$  can effectively eliminate any stickers while concurrently addressing improper stylistic issues encountered during image generation.

**Glass effect removal.** We further test  $\mathbb{W}\mathbb{M}\text{-model}$  on input images that are blemished by glass effect in Fig. 9b. We apply a fluted glass effect to images to replicate real-life scenarios where individuals capture photographs of subjects positioned behind fluted glass. This glass can have specific reflections and blurring, which may compromise the overall quality of image generation when using Textual Inversion. The use of our model can fix the distortions of the subjects and the unexpected background problem, significantly improving image quality.

## 6 Conclusion

In conclusion, we introduce ArtiFade to address the novel problem of generating high-quality and artifact-free images in the blemished subject-driven generation. Our approach involves fine-tuning a diffusion model along with artifact-free embedding to learn the alignment between unblemished images and blemished information. We present an evaluation benchmark to thoroughly assess a model’s capability in the task of blemished subject-driven generation. We demonstrate the effectiveness of ArtiFade in removing artifacts and addressing distortions in subject reconstruction under both in-distribution and out-of-distribution scenarios.

## References

- Arbel, E.; and Hel-Or, H. 2010. Shadow removal using intensity surfaces and texture anchor points. *IEEE TPAMI*, 33(6): 1202–1216.
- Avrahami, O.; Aberman, K.; Fried, O.; Cohen-Or, D.; and Lischinski, D. 2023. Break-A-Scene: Extracting Multiple Concepts from a Single Image. In *SIGGRAPH Asia 2023 Conference Papers*, 1–12.
- Bahdanau, D.; Cho, K.; and Bengio, Y. 2015. Neural machine translation by jointly learning to align and translate. In *ICLR*.
- Cao, Q.; Shen, L.; Xie, W.; Parkhi, O. M.; and Zisserman, A. 2018. Vggface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*, 67–74. IEEE.
- Cao, Z.; Niu, S.; Zhang, J.; and Wang, X. 2019. Generative adversarial networks model for visible watermark removal. *IET Image Processing*, 1783–1789.
- Caron, M.; Touvron, H.; Misra, I.; Jégou, H.; Mairal, J.; Bojanowski, P.; and Joulin, A. 2021. Emerging properties in self-supervised vision transformers. In *ICCV*, 9650–9660.
- Chen, W.; Hu, H.; Li, Y.; Ruiz, N.; Jia, X.; Chang, M.-W.; and Cohen, W. W. 2023. Subject-driven Text-to-Image Generation via Apprenticeship Learning. In Oh, A.; Neumann, T.; Globerson, A.; Saenko, K.; Hardt, M.; and Levine, S., eds., *NeurIPS*, volume 36, 30286–30305. Curran Associates, Inc.
- Chen, Z.; Long, C.; Zhang, L.; and Xiao, C. 2021. Canet: A context-aware network for shadow removal. In *ICCV*, 4743–4752.
- Chen, Z.; Zhang, Y.; Gu, J.; Zhang, Y.; Kong, L.; and Yuan, X. 2022. Cross Aggregation Transformer for Image Restoration. In *NeurIPS*.
- Cheng, D.; Li, X.; Li, W.-H.; Lu, C.; Li, F.; Zhao, H.; and Zheng, W.-S. 2018. Large-scale visible watermark detection and removal with deep convolutional networks. In *PRCV*, 27–40.
- Cheng, J.; Wu, F.; Tian, Y.; Wang, L.; and Tao, D. 2020. Rifegan: Rich feature generation for text-to-image synthesis from prior knowledge. In *CVPR*, 10911–10920.
- Dekel, T.; Rubinstein, M.; Liu, C.; and Freeman, W. T. 2017. On the effectiveness of visible watermarks. In *CVPR*, 2146–2154.
- Ding, B.; Long, C.; Zhang, L.; and Xiao, C. 2019. Argan: Attentive recurrent generative adversarial network for shadow detection and removal. In *ICCV*, 10213–10222.
- Finlayson, G.; Hordley, S.; Lu, C.; and Drew, M. 2006. On the removal of shadows from images. *IEEE TPAMI*, 28(1): 59–68.
- Finlayson, G. D.; Drew, M. S.; and Lu, C. 2009. Entropy minimization for shadow removal. *IJCV*, 85(1): 35–57.
- Finlayson, G. D.; Hordley, S. D.; and Drew, M. S. 2002. Removing shadows from images. In *ECCV*, 823–836.
- Fu, L.; Zhou, C.; Guo, Q.; Juefei-Xu, F.; Yu, H.; Feng, W.; Liu, Y.; and Wang, S. 2021. Auto-exposure fusion for single-image shadow removal. In *CVPR*, 10571–10580.
- Gal, R.; Alaluf, Y.; Atzmon, Y.; Patashnik, O.; Bermano, A. H.; Chechik, G.; and Cohen-Or, D. 2023. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *ICLR*.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative Adversarial Nets. In *NeurIPS*.
- Gu, S.; Chen, D.; Bao, J.; Wen, F.; Zhang, B.; Chen, D.; Yuan, L.; and Guo, B. 2022. Vector quantized diffusion model for text-to-image synthesis. In *CVPR*, 10696–10706.
- Guo, L.; Wang, C.; Yang, W.; Huang, S.; Wang, Y.; Pfister, H.; and Wen, B. 2023. Shadowdiffusion: When degradation prior meets diffusion model for shadow removal. In *CVPR*, 14049–14058.
- Guo, R.; Dai, Q.; and Hoiem, D. 2011. Single-image shadow detection and removal using paired regions. In *CVPR*, 2033–2040.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. In *NeurIPS*, 6840–6851.
- Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *ICLR*.
- Hu, X.; Jiang, Y.; Fu, C.-W.; and Heng, P.-A. 2019. Mask-shadowgan: Learning to remove shadows from unpaired data. In *ICCV*, 2472–2481.
- Huang, C.-H.; and Wu, J.-L. 2004. Attacking visible watermarking schemes. *IEEE TMM*, 6(1): 16–30.
- Jin, Y.; Li, R.; Yang, W.; and Tan, R. T. 2023. Estimating reflectance layer from a single image: Integrating reflectance guidance and shadow/specular aware learning. In *AAAI*, 1069–1077.
- Kawar, B.; Zada, S.; Lang, O.; Tov, O.; Chang, H.; Dekel, T.; Mosseri, I.; and Irani, M. 2023. Imagic: Text-Based Real Image Editing With Diffusion Models. In *CVPR*, 6007–6017.
- Khan, S. H.; Bennamoun, M.; Sohel, F.; and Togneri, R. 2015. Automatic shadow detection and removal from a single image. *IEEE TPAMI*, 38(3): 431–446.
- Kumari, N.; Zhang, B.; Zhang, R.; Shechtman, E.; and Zhu, J.-Y. 2023. Multi-concept customization of text-to-image diffusion. In *CVPR*, 1931–1941.
- Le, H.; and Samaras, D. 2019. Shadow removal via shadow image decomposition. In *ICCV*, 8578–8587.
- Li, B.; Qi, X.; Lukasiewicz, T.; and Torr, P. 2019a. Controllable text-to-image generation. In *NeurIPS*, volume 32.
- Li, X.; Lu, C.; Cheng, D.; Li, W.-H.; Cao, M.; Liu, B.; Ma, J.; and Zheng, W.-S. 2019b. Towards photo-realistic visible watermark removal with conditional generative adversarial networks. In *ICIG*, 345–356.
- Liang, J.; Niu, L.; Guo, F.; Long, T.; and Zhang, L. 2021. Visible watermark removal via self-calibrated localization and background refinement. In *ACM MM*, 4426–4434.

- Liu, Y.; Zhu, Z.; and Bai, X. 2021. WNet: Watermark-Decomposition Network for Visible Watermark Removal. In *WACV*, 3685–3693.
- Liu, Z.; Yin, H.; Wu, X.; Wu, Z.; Mi, Y.; and Wang, S. 2021. From shadow generation to shadow removal. In *CVPR*, 4927–4936.
- Lu, H.; Tunanyan, H.; Wang, K.; Navasardyan, S.; Wang, Z.; and Shi, H. 2023. Specialist Diffusion: Plug-and-Play Sample-Efficient Fine-Tuning of Text-to-Image Diffusion Models To Learn Any Unseen Style. In *CVPR*, 14267–14276.
- Mirza, M.; and Osindero, S. 2014. Conditional generative adversarial nets. arXiv:1411.1784.
- Mou, C.; Wang, Q.; and Zhang, J. 2022. Deep generalized unfolding networks for image restoration. In *CVPR*, 17399–17410.
- Nichol, A. Q.; Dhariwal, P.; Ramesh, A.; Shyam, P.; Mishkin, P.; McGrew, B.; Sutskever, I.; and Chen, M. 2022. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *ICML*, 16784–16804.
- Pei, S.-C.; and Zeng, Y.-C. 2006. A novel image recovery algorithm for visible watermarked images. *IEEE Trans. Inf. Forensics Secur.*, 1(4): 543–550.
- Qiao, T.; Zhang, J.; Xu, D.; and Tao, D. 2019. Mirror-gan: Learning text-to-image generation by redescription. In *CVPR*, 1505–1514.
- Qin, C.; He, Z.; Yao, H.; Cao, F.; and Gao, L. 2018. Visible watermark removal scheme based on reversible data hiding and image inpainting. *Signal Process. Image Commun.*, 60: 160–172.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*, 8748–8763.
- Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; and Chen, M. 2022. Hierarchical text-conditional image generation with clip latents. arXiv:2204.06125.
- Reed, S.; Akata, Z.; Yan, X.; Logeswaran, L.; Schiele, B.; and Lee, H. 2016. Generative adversarial text to image synthesis. In *ICML*, 1060–1069. PMLR.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *CVPR*, 10684–10695.
- Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 234–241.
- Ruan, S.; Zhang, Y.; Zhang, K.; Fan, Y.; Tang, F.; Liu, Q.; and Chen, E. 2021. Dae-gan: Dynamic aspect-aware gan for text-to-image synthesis. In *ICCV*, 13960–13969.
- Ruiz, N.; Li, Y.; Jampani, V.; Pritch, Y.; Rubinstein, M.; and Aberman, K. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 22500–22510.
- Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E. L.; Ghasemipour, K.; Gontijo Lopes, R.; Karagol Ayan, B.; Salimans, T.; et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 35: 36479–36494.
- Shor, Y.; and Lischinski, D. 2008. The shadow meets the mask: Pyramid-based shadow removal. *Comput. Graph. Forum*, 27(2): 577–586.
- Song, J.; Meng, C.; and Ermon, S. 2021. Denoising diffusion implicit models. In *ICLR*.
- Van Le, T.; Phung, H.; Nguyen, T. H.; Dao, Q.; Tran, N. N.; and Tran, A. 2023. Anti-dreambooth: Protecting users from personalized text-to-image synthesis. In *ICCV*, 2116–2127.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. In *NeurIPS*, 5998–6008.
- Wang, J.; Li, X.; and Yang, J. 2018. Stacked conditional generative adversarial networks for jointly learning shadow detection and shadow removal. In *CVPR*, 1788–1797.
- Wang, Z.; Cun, X.; Bao, J.; Zhou, W.; Liu, J.; and Li, H. 2022. Uformer: A general u-shaped transformer for image restoration. In *CVPR*, 17683–17693.
- Xiao, C.; She, R.; Xiao, D.; and Ma, K.-L. 2013a. Fast shadow removal using adaptive multi-scale illumination transfer. *Comput. Graph. Forum*, 32(8): 207–218.
- Xiao, C.; Xiao, D.; Zhang, L.; and Chen, L. 2013b. Efficient shadow removal using subregion matching illumination transfer. *Comput. Graph. Forum*, 32(7): 421–430.
- Xu, C.; Lu, Y.; and Zhou, Y. 2017. An automatic visible watermark removal technique using image inpainting algorithms. In *ICSAI*, 1152–1157.
- Xu, T.; Zhang, P.; Huang, Q.; Zhang, H.; Gan, Z.; Huang, X.; and He, X. 2018. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *CVPR*, 1316–1324.
- Yin, G.; Liu, B.; Sheng, L.; Yu, N.; Wang, X.; and Shao, J. 2019. Semantics disentangling for text-to-image generation. In *CVPR*, 2327–2336.
- Zamir, S. W.; Arora, A.; Khan, S.; Hayat, M.; Khan, F. S.; and Yang, M.-H. 2022. Restormer: Efficient transformer for high-resolution image restoration. In *CVPR*, 5728–5739.
- Zamir, S. W.; Arora, A.; Khan, S.; Hayat, M.; Khan, F. S.; Yang, M.-H.; and Shao, L. 2021. Multi-stage progressive image restoration. In *CVPR*, 14821–14831.
- Zhang, H.; Koh, J. Y.; Baldrige, J.; Lee, H.; and Yang, Y. 2021. Cross-modal contrastive learning for text-to-image generation. In *CVPR*, 833–842.
- Zhang, H.; Xu, T.; Li, H.; Zhang, S.; Wang, X.; Huang, X.; and Metaxas, D. N. 2017. Stackgan: Text to photorealistic image synthesis with stacked generative adversarial networks. In *ICCV*, 5907–5915.
- Zhang, H.; Xu, T.; Li, H.; Zhang, S.; Wang, X.; Huang, X.; and Metaxas, D. N. 2018. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE TPAMI*, 41(8): 1947–1962.
- Zhang, L.; Rao, A.; and Agrawala, M. 2023. Adding Conditional Control to Text-to-Image Diffusion Models. In *ICCV*, 3836–3847.

Zhang, L.; Zhang, Q.; and Xiao, C. 2015. Shadow Remover: Image Shadow Removal Based on Illumination Recovering Optimization. *IEEE TIP*, 24(11): 4623–4636.

Zhang, Z.; Xie, Y.; and Yang, L. 2018. Photographic text-to-image synthesis with a hierarchically-nested adversarial network. In *CVPR*, 6199–6208.

Zhu, M.; Pan, P.; Chen, W.; and Yang, Y. 2019. Dm-gan: Dynamic memory generative adversarial networks for text-to-image synthesis. In *CVPR*, 5802–5810.

Zhu, Y.; Huang, J.; Fu, X.; Zhao, F.; Sun, Q.; and Zha, Z.-J. 2022. Bijective Mapping Network for Shadow Removal. In *CVPR*, 5627–5636.

## A Training Dataset Details

Our training dataset consists of 20 training subjects, used for the fine-tuning stage of our ArtiFade models. We show an example image of each subject in Fig. 10. In Fig. 11, we showcase several unblemished images alongside their corresponding blemished versions, each featuring one of the 10 watermark types.

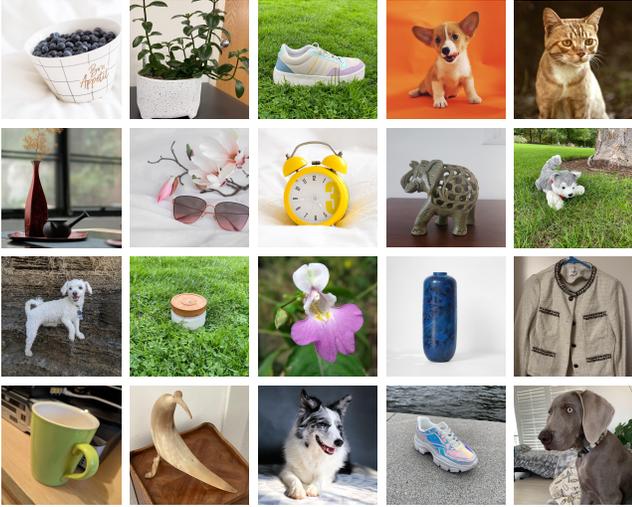


Figure 10: Examples of unblemished training images. We show a total of 20 images, each containing one distinct subject.

## B Test Dataset Details

In Fig. 12, we illustrate our `WM-ID-TEST` watermark types (see the first row) and `WM-OOD-TEST` watermark types (see the second row). The `WM-ID-TEST` watermarks are chosen from the training watermarks displayed in Fig. 11. On the other hand, the `WM-OOD-TEST` watermarks differ in font size, orientation, content, or color from all the training watermarks presented in Fig. 11.

## C Analysis of Watermark Density

In Fig. 13, we present results to illustrate the impact of varying watermark densities (*i.e.*, varying qualities), highlighting the robust ability of our `WM-model` to remove watermarks under all conditions.

## D Analysis of Unblemished Image Ratio

We employ our `WM-model` to evaluate the performance when the input images contain different proportions of unblemished images. We test our `WM-model` and Textual Inversion on five ratios of unblemished images: 100%, 75%, 50%, 25%, and 0%. The results are shown in Fig. 14.

Notably, even when there is only one blemished image in the second column example, the impact on Textual Inversion is already evident, which deteriorates as the ratio decreases. Instead, our method effectively eliminates artifacts in all settings of unblemished image ratio, demonstrating its versatility in real-life scenarios.

## E Analysis of Training Dataset Size

We conduct an analysis to investigate the impact of the number of training subjects (*i.e.*, the size of the training dataset) on the performance of our model. We utilize the same set of artifacts  $L_{WM} = 10$ , as described in Method in the main paper. We construct blemished training datasets in four different sizes: (1) with 5 subjects, (2) with 10 subjects, (3) with 15 subjects, and (4) with 20 subjects. We generate 50, 100, 150, and 200 blemished datasets for each of these cases. Subsequently, we fine-tune four distinct ArtiFade models, each with 16k training steps.

We compare the models trained using different data sizes under the in-distribution scenario (see Fig. 15a) and under the out-of-distribution scenario (see Fig. 15b). We note that when the number of training subjects is less than 15,  $I^{DINO}$  and  $T^{CLIP}$  are relatively lower than the other two cases in both ID and OOD scenarios. This observation can be attributed to a significant likelihood of subject or background overfitting during the reconstruction and image synthesis processes, as visually illustrated in Fig. 16 and Fig. 17. However, as the number of training subjects reaches or exceeds 15, we observe a convergence in the values of  $I^{DINO}$  and  $T^{CLIP}$ , indicating a reduction in subject overfitting. Regarding  $R^{DINO}$ , we note that all cases exhibit values greater than one, with a slightly increasing trend as the number of training subjects rises.

## F Failure Cases

We present several failure cases when applying ArtiFade based on Textual Inversion. We demonstrate the limitations of our `WM-model` in Fig. 18. Despite the model’s ability to eliminate watermarks, we still encounter issues with incorrect subject color, as shown in Fig. 18a, which arises due to the influence of the watermark color. We also encounter incorrect subject identity in some cases, as demonstrated in Fig. 18b. One possible reason is that the watermarks significantly contaminate the images, causing the learning process of embedding to focus on the contaminated visual appearance instead of the intact subject. Another failure case is subject overfitting, as shown in Fig. 18c. In this case, the constructed subject overfits with a similar subject type that appears in the training dataset. This problem occurs because the blemished embedding of the testing subject closely resembles some blemished embeddings of the training subjects. Surprisingly, we find those problems can be solved by using ArtiFade based on DreamBooth, which is mentioned in Sec. 4.5. Therefore, we recommend using ArtiFade based on DreamBooth when encountering the limitations mentioned above.

## G Additional Comparison with Textual Inversion

We use the same training subjects with  $N = 20$  from Sec. 3.3 to train an ArtiFade model named `RC-model` using red circle artifacts. For the training set of `RC-model`, due to the simplicity of red circles, we only synthesize a single blemished subset (*i.e.*,  $L_{RC} = 1$ ) for each subject, deriving 20 blemished subsets in total. We augment each image with a

Unblemished image

Blemished images



Figure 11: Examples of the training dataset: unblemished images and their corresponding blemished images.



Figure 12: Example of test watermark types. The first row displays the WM-ID-TEST, while the second row presents the WM-OOD-TEST.

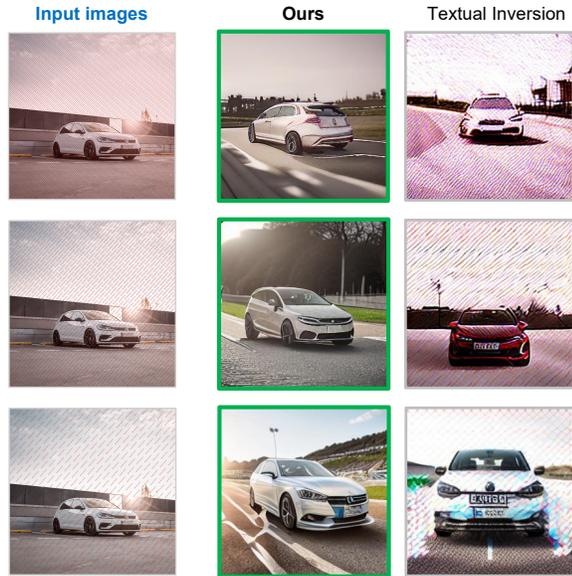


Figure 13: Varying qualities of input images. Our method ( $\mathcal{W}\mathcal{M}\text{-model}$ ) can be used to remove watermarks when input images are of any quality.

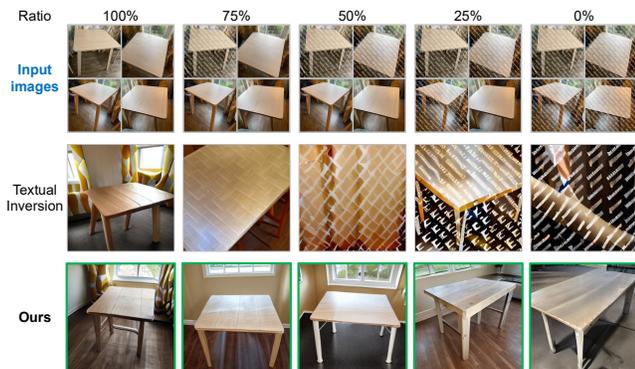


Figure 14: Comparison between different ratios of unblemished images. ArtiFade can perform well under any scenarios with different ratios of unblemished images.

red circle mark that is randomly scaled and positioned on the source image. Considering the small scale of  $\mathcal{R}\mathcal{C}\text{-model}$ 's

Method	RC-test				
	I <sup>DINO</sup>	R <sup>DINO</sup>	I <sup>CLIP</sup>	R <sup>CLIP</sup>	T <sup>CLIP</sup>
TI (unblemished)	0.488	1.021	0.730	1.077	0.283
TI (blemished)	0.406	0.990	0.672	1.042	0.284
Ours ( $\mathcal{R}\mathcal{C}\text{-model}$ )	<b>0.476</b>	<b>1.013</b>	0.722	<b>1.065</b>	<b>0.285</b>
Ours ( $\mathcal{W}\mathcal{M}\text{-model}$ )	0.474	1.006	<b>0.727</b>	1.063	0.282

Table 5: Quantitative results of RC-test.

datasets, we only fine-tune  $\mathcal{R}\mathcal{C}\text{-model}$  for 8k steps. We further introduce RC-test, which applies only one type of artifact (*i.e.*, red circle) to our 16 test subjects, resulting in 16 test sets. We test both  $\mathcal{R}\mathcal{C}\text{-model}$  and  $\mathcal{W}\mathcal{M}\text{-model}$  on RC-test. The quantitative and qualitative results are shown in Tab. 5 and Fig. 19, respectively.

**Quantitative results analysis.** From Tab. 5, we can observe that both  $\mathcal{R}\mathcal{C}\text{-model}$  and  $\mathcal{W}\mathcal{M}\text{-model}$  yield higher results in nearly all cases than Textual Inversion (Gal et al. 2023) with blemished inputs, showing the capability of our models to eliminate artifacts and generate subjects with higher fidelity. It is important to note that the RC-test is considered out-of-distribution with respect to  $\mathcal{W}\mathcal{M}\text{-model}$ . Nevertheless, the metrics produced by  $\mathcal{W}\mathcal{M}\text{-model}$  remain comparable to those of  $\mathcal{R}\mathcal{C}\text{-model}$ , with a minor difference observed. These results provide additional evidence supporting the generalizability of our  $\mathcal{W}\mathcal{M}\text{-model}$ .

**Qualitative results analysis.** As illustrated in Fig. 19, Textual Inversion struggles with accurate color reconstruction. It also showcases subject distortions and introduces red-circle-like artifacts during image generation when using blemished embeddings. In contrast, our  $\mathcal{R}\mathcal{C}\text{-model}$  (see Fig. 19a) and  $\mathcal{W}\mathcal{M}\text{-model}$  (see Fig. 19b) are capable of generating high-quality images that accurately reconstruct the color and identities of subjects without any interference from artifacts during the image synthesis.

## H Additional Qualitative Comparisons

We present additional qualitative results comparing our ArtiFade models with Textual Inversion (Gal et al. 2023) and DreamBooth (Ruiz et al. 2023) in Fig. 20. We employ  $\mathcal{W}\mathcal{M}\text{-model}$  and ArtiFade based on DreamBooth mentioned in Sec. 4.5. Textual Inversion generates images with distorted subjects and backgrounds contaminated by watermarks, whereas DreamBooth can effectively capture intricate subject details and accurately reproduce watermark patterns. In contrast, our models (*i.e.*, TI-based and DB-based ArtiFade) generate images devoid of watermark pollution with correct subject identities for both in-distribution (see the first three rows in Fig. 20) and out-of-distribution (see the last two rows in Fig. 20) cases. Notably, our method based on DreamBooth preserves the high fidelity and finer detail reconstruction benefits of vanilla DreamBooth, even in the context of blemished subject-driven generation.

In Fig. 21, we show qualitative results for subjects with complex features (*e.g.*, human faces) using our models, Textual Inversion, DreamBooth and Break-a-Scene (Avrahami

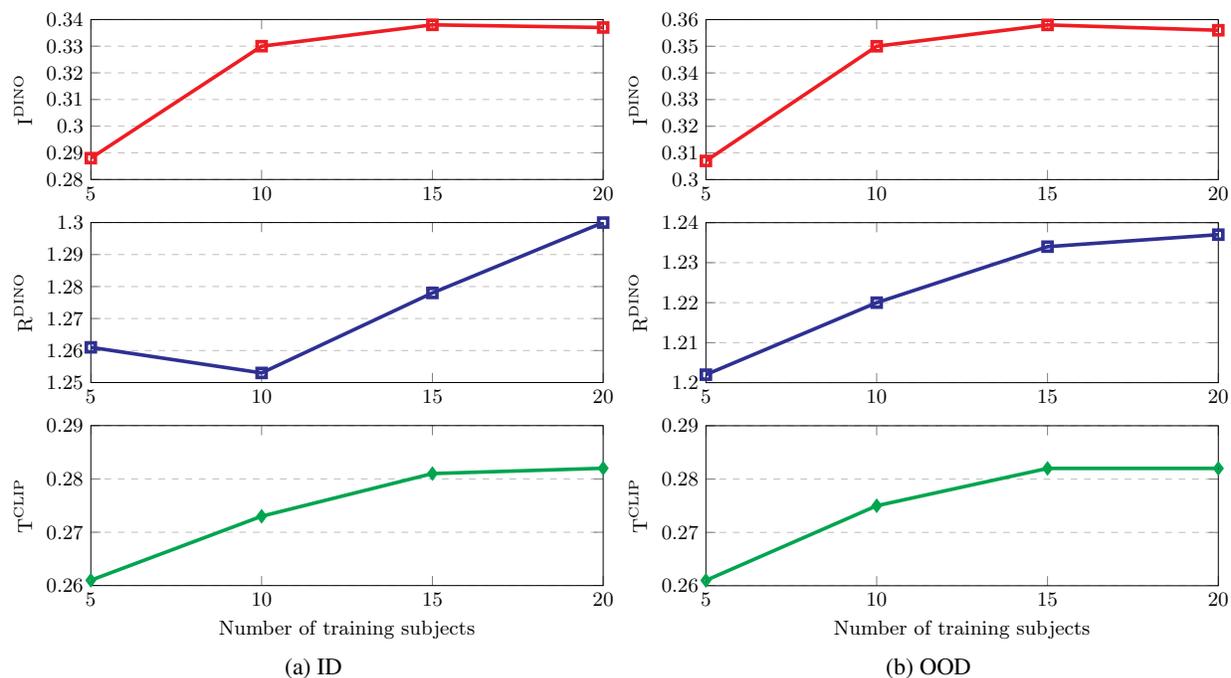


Figure 15: Analysis of the number of training subjects.

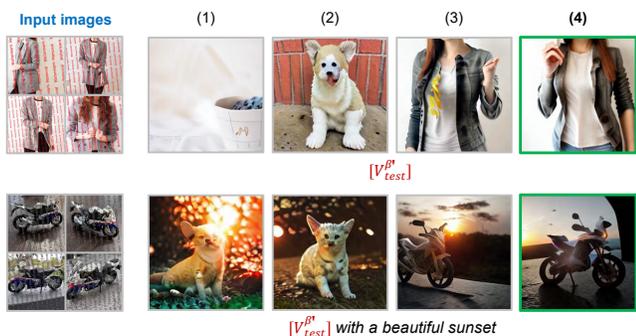


Figure 16: Qualitative results of different number of training subjects - ID.

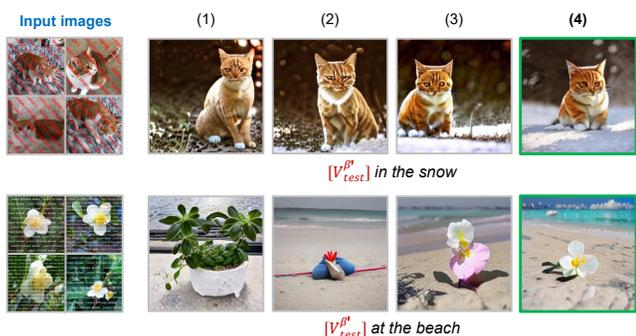


Figure 17: Qualitative results of different number of training subjects - OOD.

et al. 2023). Break-a-Scene can separate multiple subjects inside one image. We use Break-a-scene to generate human-

only images. However, we find that Break-a-scene fails to separate humans from artifacts, resulting in polluted images. As a result, our methods (*i.e.*, TI-based and DB-based Arti-Fade) consistently surpass Textual Inversion, DreamBooth, and Break-a-Scene, achieving high-quality image generation of complex data in in-distribution cases, as shown in the first two rows of Fig. 21, and out-of-distribution cases, as illustrated in the last row of Fig. 21.

## I More Applications

We explore more applications of our `WM_model`, demonstrating its versatility beyond watermark removal. As shown in Fig. 22, our model exhibits the capability to effectively eliminate unwanted artifacts from images, enhancing their visual quality. Furthermore, our model showcases the ability to recover incorrect image styles induced by artifacts, thereby restoring the intended style of the images.

## J Social Impact

Our research addresses the emerging challenge of generating content from images with embedded watermarks, a scenario we term blemished subject-driven generation. Users often source images from the internet, some of which may contain watermarks intended to protect the original author’s copyright and identity. However, our method is capable of removing various types of watermarks, potentially compromising the authorship and copyright protection. This could lead to increased instances of image piracy and the generation of illicit content. Hence, we advocate for legal compliance and the implementation of usage restrictions to govern the deployment of our technique and subsequent models in the future.

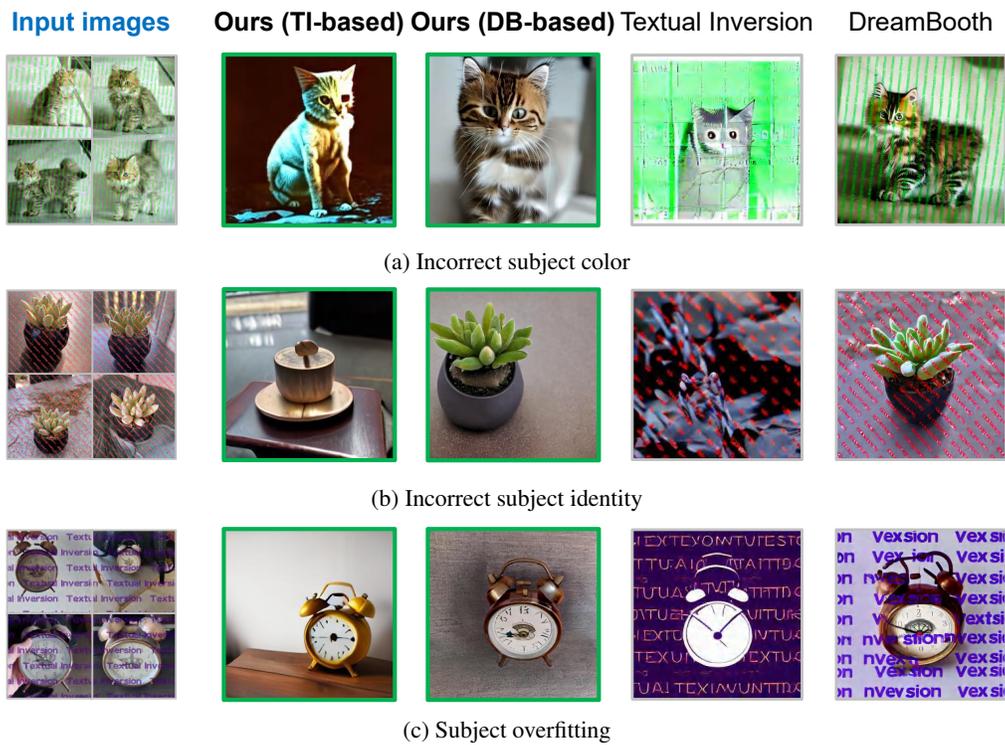
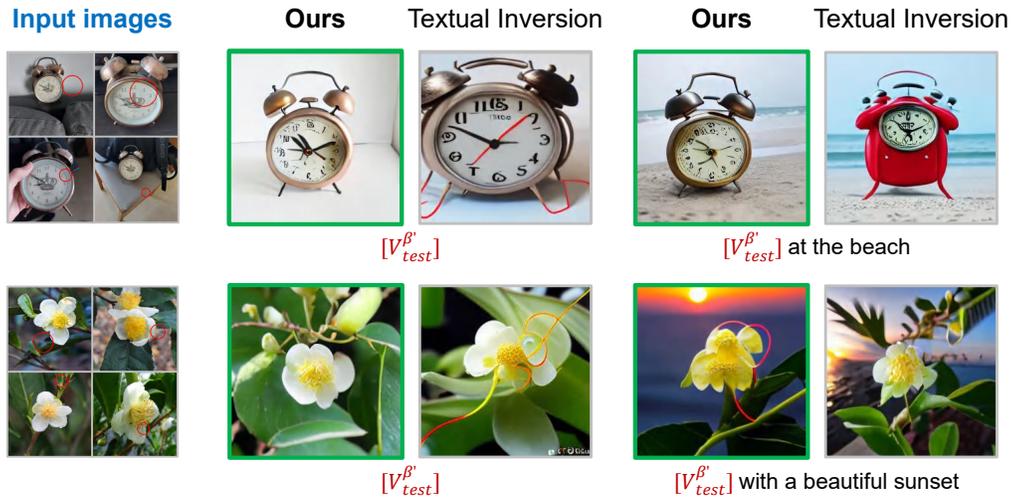
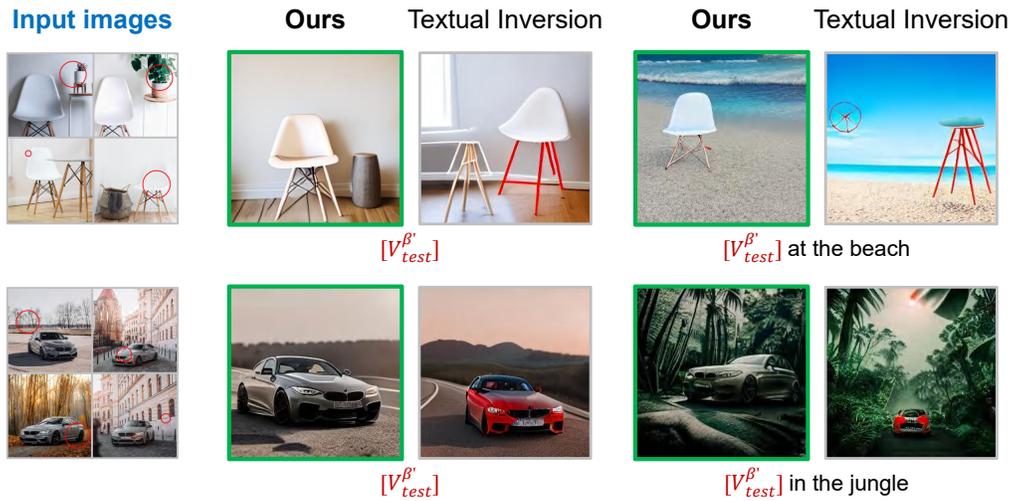


Figure 18: Failure cases of ArtiFade based on Textual Inversion. We observe three main types of failure cases of our `WM_model`: (a) incorrect subject color, (b) incorrect subject identity, and (c) subject overfitting. However, those limitations can be resolved by using ArtiFade with DreamBooth-based fine-tuning.



(a) RC-model on RC-test.



(b) WM-model on RC-test.

Figure 19: Qualitative results of RC-test. Our models consistently output high-quality and artifact-free images compared to Textual Inversion.

Input images

Ours (TI-based) Ours (DB-based) Textual Inversion

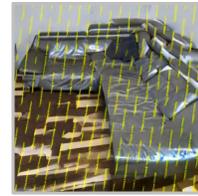
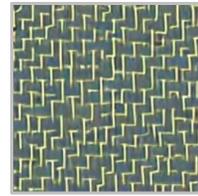
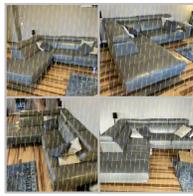
DreamBooth



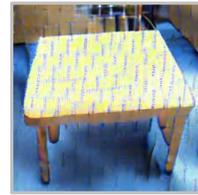
$[V^{\beta'}_{test}]$



$[V^{\beta'}_{test}]$



$[V^{\beta'}_{test}]$



$[V^{\beta'}_{test}]$



$[V^{\beta'}_{test}]$

Figure 20: Additional qualitative comparisons.

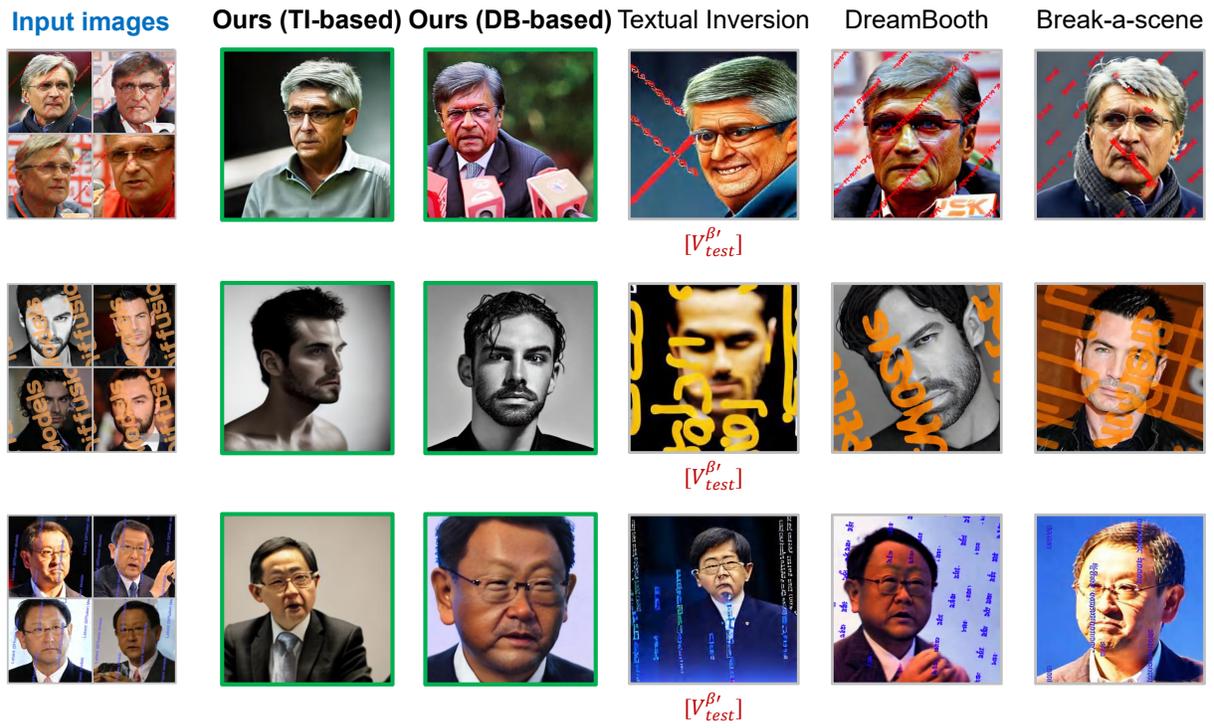


Figure 21: Additional qualitative comparisons - Human Faces.

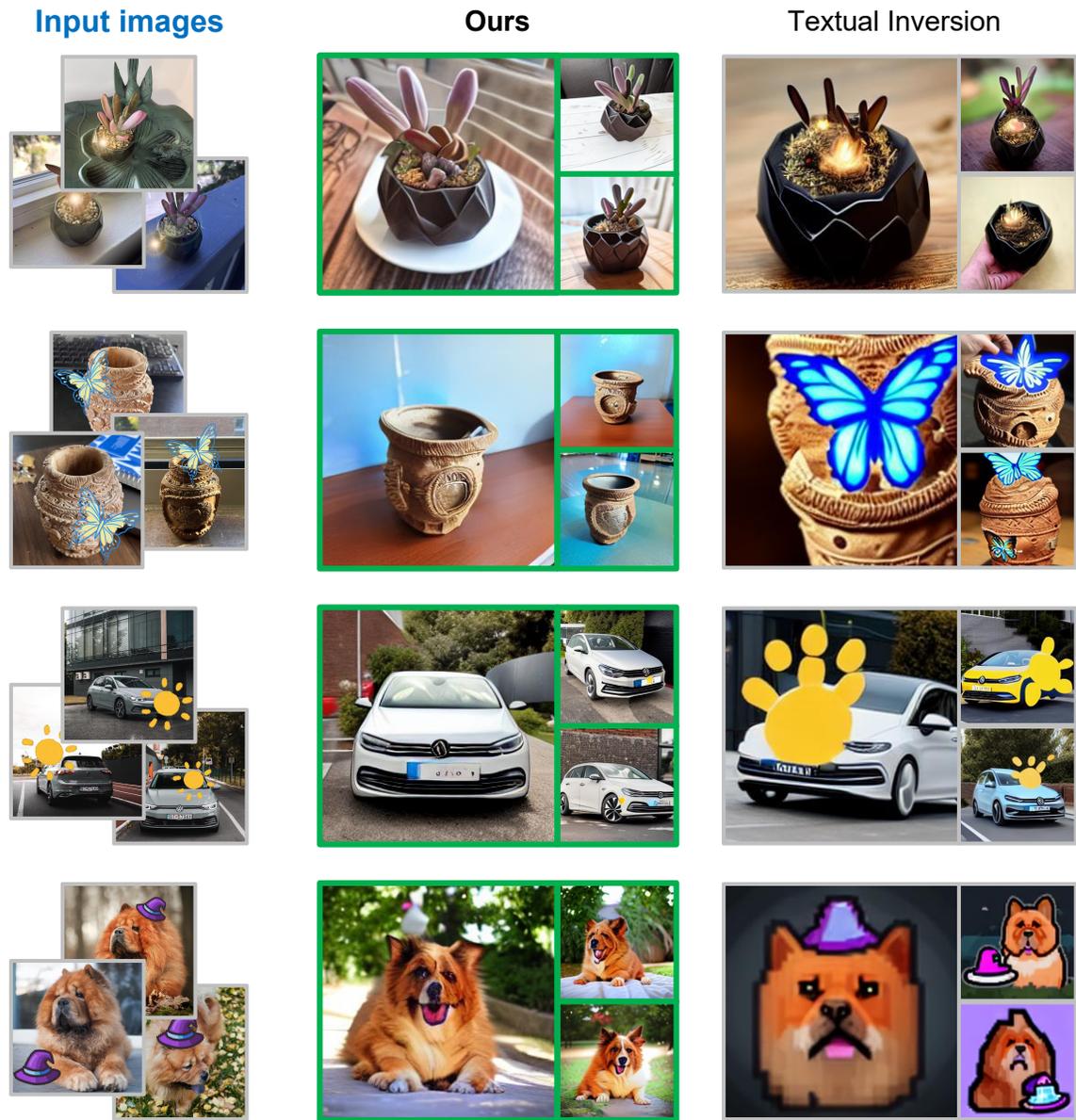


Figure 22: More applications. Our `WM-model` can be used to eliminate various stickers and fix the incorrect image style.