# Anti-DreamBooth: Protecting users from personalized text-to-image synthesis

Thanh Van Le, Hao Phung, Thuan Hoang Nguyen, Quan Dao, Ngoc Tran, Anh Tran
VinAI Research

{thanhlv19, haopt12, thuannh5, quandm7, ngoctnq, anhtt152}@vinai.io

## Abstract

*Text-to-image diffusion models are nothing but a revolution, allowing anyone, even without design skills, to create realistic images from simple text inputs. With powerful personalization tools like DreamBooth, they can generate images of a specific person just by learning from his/her few reference images. However, when misused, such a powerful and convenient tool can produce fake news or disturbing content targeting any individual victim, posing a severe negative social impact. In this paper, we explore a defense system called Anti-DreamBooth against such malicious use of DreamBooth. The system aims to add subtle noise perturbation to each user's image before publishing in order to disrupt the generation quality of any DreamBooth model trained on these perturbed images. We investigate a wide range of algorithms for perturbation optimization and extensively evaluate them on two facial datasets over various text-to-image model versions. Despite the complicated formulation of DreamBooth and Diffusion-based text-to-image models, our methods effectively defend users from the malicious use of those models. Their effectiveness withstands even adverse conditions, such as model or prompt/term mismatching between training and testing. Our code will be available at* https://github.com/VinAIResearch/Anti-DreamBooth.git.

## 1. Introduction

Within a few years, denoising diffusion models [23, 51, 41] have revolutionized image generation studies, allowing producing images with realistic quality and diverse content [19]. They especially succeed when being combined with language [38] or vision-language models [37] for text-to-image generation. Large models [39, 3, 41, 45, 9] can produce photo-realistic or artistic images just from simple text description inputs. A user can now generate art within a few seconds, and a generated drawing even beat professional artists in an art competition [42]. Photo-realistic synthetic images can be hard to distinguish from real photos [26]. Besides, ControlNet [59] offers extra options to control the
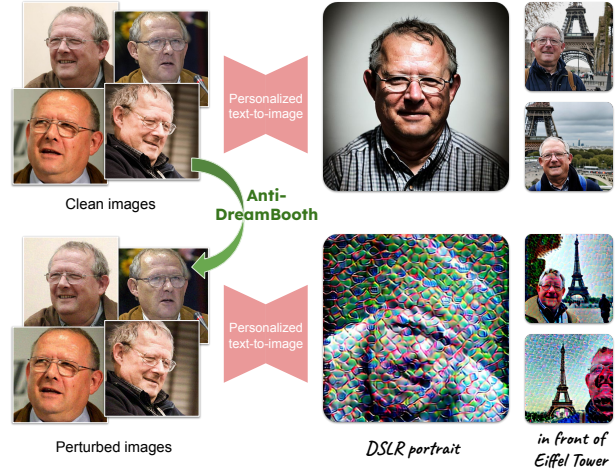


Figure 1: A malicious attacker can collect a user's images to train a personalized text-to-image generator for malicious purposes. Our system, called Anti-DreamBooth, applies imperceptible perturbations to the user's images before releasing, making any personalized generator trained on these images fail to produce usable images, protecting the user from that threat.

generation outputs, further boosting the power of the text-to-image models and bringing them closer to mass users.

One extremely useful feature for image generation models is personalization, which allows the models to generate images of a specific subject, given a few reference examples. For instance, one can create images of himself/herself in a fantasy world for fun, or create images of his/her family members as a gift. Textual Inversion [20] and DreamBooth [44] are two prominent techniques that offer that impressive ability. While Textual Inversion only optimizes the text embedding inputs representing the target subject, DreamBooth finetunes the text-to-image model itself for better personalization quality. Hence, DreamBooth is particularly popular and has become the core technique in many applications.

While the mentioned techniques provide a powerful and convenient tool for producing desirable images at will, they also pose a severe risk of being misused. A malicious user can propagate fake news with photo-realistic images of a

celebrity generated by DreamBooth. This can be classified as DeepFakes [27], one of the most serious AI crime threats that has drawn an enormous attention from the media and community in recent years. Besides creating fake news, DreamBooth can be used to issue harmful images targeting specific persons, disrupting their lives and reputations. While the threat of GAN-based DeepFakes techniques is well-known and has drawn much research interest, the danger from DreamBooth has yet to be aware by the community, making its damage, when happening, more dreadful.

This paper discusses how to protect users from malicious personalized text-to-image synthesis. Inspired by Deep-Fakes's prevention studies [57, 43, 56, 25, 55], we propose to pro-actively defend each user from the DreamBooth threat by injecting subtle adversarial noise into their images before publishing. The noise is designed so that any DreamBooth model trained on these perturbed images fails to produce reasonable-quality images of the target subject. While the proposed mechanism, called Anti-DreamBooth, shares the same goal and objective as the techniques to disrupt GAN-based DeepFakes, it has a different nature due to the complex formulation of diffusion-based text-to-image models and DreamBooth:

- In GAN-based disruption techniques, the defender optimizes the adversarial noise of a single image, targeting a fixed DeepFakes generator. In Anti-DreamBooth, we have to optimize the perturbation noise to disrupt a dynamic, unknown generator that is finetuned from the perturbed images themselves.

- GAN-based DeepFakes generator produces each fake image via a single forward step; hence, adversarial noise can be easily learned based on the model's gradient. In contrast, a diffusion-based generator produces each output image via a series of non-deterministic denoising steps, making it impossible to compute the end-to-end gradient for optimization.

- Anti-DreamBooth has a more complex setting by considering many distinctive factors, such as the prompt used in training and inference, the text-to-image model structure and pre-trained weights, and more.

Despite the complexity mentioned above, we show that the DreamBooth threat can be effectively prevented. Instead of targeting the end-to-end image generation process, we can adapt the adversarial learning process to break each diffusion sampling step. We design different algorithms for adversarial noise generation, and verify their effectiveness in defending DreamBooth attack on two facial benchmarks. Our proposed algorithms successfully break all DreamBooth attempts in the controlled settings, causing the generated images to have prominent visual artifacts. Our proposed defense shows consistent effect when using different text-to-image models and different training text prompts. More impressively, Anti-DreamBooth maintains its efficiency even under adverse conditions, such as model or prompt/term mismatching between training and testing.

In summary, our contributions include: (1) We discuss the potential negative impact of personalized text-to-image synthesis, particularly with DreamBooth, and define a new task of defending users from this critical risk, (2) We propose proactively protecting users from the threat by adding adversarial noise to their images before publishing, (3) We design different algorithms for adversarial noise generation, adapting to the step-based diffusion process and finetuning-based DreamBooth procedure, (4) We extensively evaluate our proposed methods on two facial benchmarks and under different configurations. Our best defense works effectively in both convenient and adverse settings.

## 2. Related work

### 2.1. Text-to-image generation models

Due to the advent of new large-scale training datasets such as LAION5B [47], text-to-image generative models are advancing rapidly, opening new doors in many visual-based applications and attracting attention from the public. These models can be grouped into four main categories: auto-regressive [58], mask-prediction [13], GAN-based [46] and diffusion-based approaches, all of which show astounding qualitative and quantitative results. Among these methods, diffusion-based models [41, 45, 9, 36, 40] have exhibited an exceptional capacity for generating high-quality and easily modifiable images, leading to their widespread adoption in text-to-image synthesis. GLIDE [36] is arguably the first to combine a diffusion model with classifier guidance for text-to-image generation. DALL-E 2 [40] then improves the quality further using the CLIP text encoder and diffusion-based prior. For better trade-off between efficiency and fidelity, following-up works either introduce coarse-to-fine generation process like Imagen [45] and eDiff-I [9] or work on latent space like LDM [41]. StableDiffusion [4], primarily based on LDM, is the first open-source large model of this type, further boosting the widespread applications of text-to-image synthesis.

### 2.2. Personalization

Customizing the model's outputs for a particular person or object has been a significant aim in the machine-learning community for a long time. Generally, personalized models are commonly observed in recommendation systems [5] or federated learning [48, 22]. Within the context of diffusion models, previous research has focused on adapting a pre-trained model to create fresh images based on a particular target idea using natural language cues. Existing methods for personalizing the model either involve adjusting a col-

lection of text embeddings to describe the idea [20] or fine-tuning the denoising network to connect a less commonly used word-embedding to the novel concept [44]. For better customization, [29] propose a novel approach to not only model a new concept of an individual subject by optimizing a small set of parameters of cross-attention layers but also to combine multiple concepts of objects via joint training. Among these tools, DreamBooth [44] is particularly popular due to its exceptional quality and has become the core technique in many applications. Hence, we focus on defending the risks of malign image generation coming from this technique.

## 2.3. Adversarial attacks

With the introduction of the Fast Gradient Sign Method (FGSM) attack [21], adversarial vulnerability has become an active field of research in machine learning. The goal of adversarial attacks is to generate a model input that can induce a misclassification while remaining visually indistinguishable from a clean one. Following this foundational work, different attacks with different approaches started to emerge, with more notable ones including: [30, 33] being FGSM's iterative versions, [12] limiting the adversarial perturbation's magnitude implicitly using regularization instead of projection, [35] searching for a close-by decision boundary to cross, etc. For black-box attacks, where the adversary does not have full access to the model weights and gradients, [53] estimates the gradient using sampling methods, while [10, 14, 6] aim to synthesize a close-by example by searching for the classification boundary, then finding the direction to traverse towards a good adversarial example. Combining them, [16] is an ensemble of various attacks that are commonly used as a benchmark metric, being able to break through gradient obfuscation [7] with the expectation-over-transformation technique [8].

## 2.4. User protection with image cloaking

With the rapid development of AI models, their misuse risk has emerged and become critical. Particularly, many models exploit the public images of each individual for malicious purposes. Instead of passively detecting and mitigating these malign actions, many studies propose proactively preventing them from succeeding. The idea is to add subtle noise into users' images before publishing to disrupt any attempt to exploit those images. This approach is called "image cloaking", which our proposed methods belong to.

One application of image cloaking is to prevent privacy violations caused by unauthorized face recognition systems. Fawkes [49] applies targeted attacks to shift the user's identity towards a different reference person in the embedding space. Although it learns the adversarial noise using a surrogate face recognition model, the noise successfully transfers to break other black-box recognizers. Lowkey [15] fur-

ther improves the transferability by using an ensemble of surrogate models. It also considers a Gaussian smoothed version of the perturbed image in optimization, improving robustness against different image transformations. AMT-GAN [24] crafts a natural-looking cloak via makeup transfer, while OPOM [60] optimizes person-specific universal privacy masks.

Another important application of image cloaking is to disrupt GAN-based image manipulation for DeepFakes. Yang et al. [56] exploits differentiable image transformations for robust image cloaking. Yeh et al. [57] defines new effective objective functions to nullify or distort the image manipulation. Huang et al. [25] addresses personalized DeepFakes techniques by alternating the training of the surrogate model and a perturbation generator. Anti-Forgery [55] crafts the perturbation for channels $a$ and $b$ in the $Lab$ color space, aiming for natural-looking and robust cloaking. Lately, UnGANable [31] prevents StyleGAN-based image manipulation by breaking its inversion process.

## 3. Problem

### 3.1. Background

**Adversarial attacks.** The goal of adversarial attacks is to find an imperceptible perturbation of an input image to mislead the behavior of given models. Typical works have been developed for classification problems where for a model $f$, an adversarial example $x'$ of an input image $x$ is generated to stay visually undetectable while inducing a misclassification $y_{\text{true}} \neq f(x')$ (untargeted attack), or making the model predict a predefined target label $y_{\text{target}} = f(x_{\text{adv}}) \neq y_{\text{true}}$ (targeted attack). The minimal visual difference is usually enforced by bounding the perturbation to be within an $\eta$-ball w.r.t. an $\ell_p$ metrics, that is $\|x' - x\|_p < \eta$. To achieve this objective, denoting $\Delta = \{\delta : \|\delta\|_p \leq \eta\}$, we find the optimal perturbation $\delta$ to maximize the classification loss in the untargeted version:

$$\delta_{\text{adv}} = \arg\max_{\delta \in \Delta} \mathcal{L}(f(x + \delta), y_{\text{true}}), \quad (1)$$

or to minimize the loss for the targeted variant:

$$\delta_{\text{adv}} = \arg\min_{\delta \in \Delta} \mathcal{L}(f(x + \delta), y_{\text{target}}). \quad (2)$$

Projected Gradient Descent (PGD) [33] is a commonly used attack based on an iterative optimization process. The updating pipeline to predict $x'$ for untargeted attack is:

$$\begin{aligned} x'_0 &= x \\ x'_k &= \Pi_{(x,\eta)}(x'_{k-1} + \alpha \cdot \text{sgn}(\nabla_x \mathcal{L}(f(x'_{k-1}), y_{true}))) \end{aligned} \quad (3)$$

where $\Pi_{x,\eta}(z)$ restrains pixel values of $z$ within an $\eta$-ball around the original values in $x$. We acquire the adversarial example $x'$ after a pre-defined number of iterations.

**Diffusion models** are a type of generative models [50, 23] that decouple the role of generation into two opposing procedures: a forward process and a backward process. While the forward process gradually adds noise to an input image until data distribution becomes pure Gaussian noise, the latter learns to reverse that process to obtain the desired data from random noise. Given input image $x_0 \sim q(x)$, the diffusion process perturbs the data distribution with a noise scheduler $\{\beta_t : \beta_t \in (0, 1)\}_{t=1}^T$ producing increasing levels of noise addition through $T$ steps to obtain a sequence of noisy variables: $\{x_1, x_2, \ldots, x_T\}$. Each variable $x_t$ is constructed via injecting noise at corresponding timestep $t$:

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon \qquad (4)$$

where $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ and $\epsilon \sim \mathcal{N}(0, \mathbf{I})$.

The backward process learns to denoise from noisy variable $x_{t+1}$ to less-noisy variable $x_t$ via simply estimating the injected noise $\epsilon$ with a parametric neural network $\epsilon_\theta(x_{t+1}, t)$. The denoising process is trained to minimize $\ell_2$ distance between estimated noise and true noise:

$$\mathcal{L}_{unc}(\theta, x_0) = \mathbb{E}_{x_0, t, \epsilon \in \mathcal{N}(0,1)} \|\epsilon - \epsilon_\theta(x_{t+1}, t)\|_2^2 \qquad (5)$$

where $t$ is uniformly samples within $\{1, \ldots, T\}$.

**Prompt-based Diffusion Models.** Unlike unconditional diffusion models, prompt-based diffusion models control the sampling process with an additional prompt $c$ to generate photo-realistic outputs which are well-aligned with the text description. The objective is formulated as follows:

$$\mathcal{L}_{cond}(\theta, x_0) = \mathbb{E}_{x_0, c, \epsilon \in \mathcal{N}(0,1)} \|\epsilon - \epsilon_\theta(x_{t+1}, t, c)\|_2^2 \qquad (6)$$

Thanks to prompt condition, the model can produce more excellent performance in terms of visual quality than its unconditional counterparts. However, the implementation of most prominent methods [45, 9] are not publicly available. Alternatively, Stable Diffusion has released the pre-trained weights based on Hugging Face implementation [54] to facilitate research in the community. Hence, we mainly perform experiments on different versions of Stable Diffusion. **DreamBooth** is a finetuning technique to personalize text-to-image diffusion models for instance of interest. This technique has two aims. First, it enforces the model to learn to reconstruct the user's images, with a generic prompt $c$ such as "a photo of *sks* [class noun]", with *sks* is a special term denoting the target user, and "[class noun]" is the object type, which can be "person" for human subject. To train this, DreamBooth employs the base loss of diffusion models in Eq. (6), with $x_0$ is each user's reference image. Second, it further introduces a prior preservation loss to prevent overfitting and text-shifting problems when only a small set of instance examples is used. More precisely, it uses a generic prior prompt $c_{pr}$, e.g.,"a photo of [class noun]", and enforces the model to reproduce instance examples randomly

generated from that prior prompt using the original weights $\theta_{ori}$. The training objective is the combination of two objectives:

$$\mathcal{L}_{db}(\theta, x_0) = \mathbb{E}_{x_0, t, t'} \|\epsilon - \epsilon_\theta(x_{t+1}, t, c)\|_2^2$$
$$+ \lambda \|\epsilon' - \epsilon_{\theta_{ori}}(x_{t'+1}, t', c_{pr})\|_2^2 \qquad (7)$$

where $\epsilon, \epsilon'$ are both sampled from $\mathcal{N}(0, \mathbf{I})$ and $\lambda$ emphasizes the importance of the prior term. While DreamBooth was originally designed for Imagen [45], it was quickly adopted for any text-to-image generator.

### 3.2. Problem definition

DreamBooth is a powerful tool to generate photo-realistic outputs of a target instance with rich context beyond the range of reference samples. With such an impressive capacity, DreamBooth can be a double-edged sword. When misused, it can generate harmful images toward the target individual. To mitigate this phenomenon, we propose to craft an imperceptible perturbation added to each user's image that can disrupt the finetuned DreamBooth models to generate distorted images with noticeable artifacts. We define the problem formally below.

Denote $\mathcal{X}$ as the set of images of the person to protect. For each image $x \in \mathcal{X}$, we add an adversarial perturbation $\delta$ and publish the modified image $x' = x + \delta$, while keeping the original one private. The published image set is called $\mathcal{X}'$. An adversary can collect a small image set of that person $\mathcal{X}'_{db} = \{x^{(i)} + \delta^{(i)}\}_{i=1}^{N_{db}} \subset \mathcal{X}'$. He then uses that set as a reference to finetune a text-to-image generator $\epsilon_\theta$, following the DreamBooth algorithm, to get the optimal hyper-parameters $\theta^*$. The general objective is to optimize the adversarial noise $\Delta_{db} = \{\delta^{(i)}\}_{i=1}^{N_{db}}$ that minimizes the personalized generation ability of that DreamBooth model:

$$\Delta_{db}^* = \underset{\Delta_{db}}{\arg \min} \, \mathcal{A}(\epsilon_{\theta^*}, \mathcal{X}),$$

$$\text{s.t.} \quad \theta^* = \underset{\theta}{\arg \min} \sum_{i=1}^{N_{db}} \mathcal{L}_{db}(\theta, x^{(i)} + \delta^{(i)}), \qquad (8)$$

$$\text{and} \quad \|\delta^{(i)}\|_p \leq \eta \quad \forall i \in \{1, 2, .., N_{db}\},$$

where $\mathcal{L}_{db}$ is defined in Eq. 7 and $\mathcal{A}(\epsilon_{\theta^*}, \mathcal{X})$ is some personalization evaluation function that assesses the quality of images generated by the DreamBooth model $\epsilon_{\theta^*}$ and the identity correctness based on the reference image set $\mathcal{X}$.

However, it is hard to define a unified evaluation function $\mathcal{A}$. A defense succeeds when the DreamBooth-generated images satisfy one of the criteria: (1) awful quality due to extreme noise, blur, distortion, or noticeable artifacts, (2) shifted content with none or unrecognizable human subjects, (3) mismatched subject identity. Even when we focus on the first criteria, there is no all-in-one image quality assessment metric. Instead, we can use simpler objective

functions disrupting the DreamBooth training to achieve the same goal.

We further divide the defense settings into categories, from easy to hard: convenient, adverse, and uncontrolled.

**Convenient setting.** In this setting, we have prior knowledge about the pretrained text-to-image generator, training term (e.g., "sks"), and training prompt $c$ the attacker will use. While this setting sounds restricted, it is practical. First, the pretrained generator has to be high-quality and open-source. So far, only Stable Diffusion has been made publicly available with several versions released. Second, people often use the default training term and prompt provided in DreamBooth's code. This setting can be considered as "white-box".

**Adverse settings.** In these settings, the pretrained text-to-image generator, training term, or training prompt used by the adversary is unknown. The defense method, if needed, can use a surrogate component that potentially mismatches the actual one to craft the adversarial noises. These settings can be considered as "gray-box".

**Uncontrolled setting.** This is an extra, advanced setting in which some of the user's clean images are leaked to the public without our control. The adversary, therefore, can collect a mix of perturbed and clean images $\mathcal{X}'_{db} = \mathcal{X}'_{adv} \cup \mathcal{X}_{cl}$, with $\mathcal{X}'_{adv} \subset \mathcal{X}'$ and $\mathcal{X}_{cl} \subset \mathcal{X}$. This setting is pretty challenging since the DreamBooth model can learn from unperturbed photos to generate reasonable personalized images.

# 4. Proposed defense methods

## 4.1. Overall direction

As discussed, instead of defining some evaluation function $\mathcal{A}$ for optimization (Eq. 8), we can aim to attack the learning process of DreamBooth. As the DreamBooth model overfits the adversarial images, we can trick it into performing worse in reconstructing clean images:

$$\delta^{*(i)} = \arg\max_{\delta^{(i)}} \mathcal{L}_{cond}(\theta^*, x^{(i)}), \forall i \in \{1,..,N_{db}\},$$

$$\text{s.t.} \quad \theta^* = \arg\min_\theta \sum_{i=1}^{N_{db}} \mathcal{L}_{db}(\theta, x^{(i)} + \delta^{(i)}), \qquad (9)$$

$$\text{and} \quad \|\delta^{(i)}\|_p \leq \eta \quad \forall i \in \{1,..,N_{db}\},$$

where $\mathcal{L}_{cond}$ and $\mathcal{L}_{db}$ are defined in Eq. 6 and 7. Note that, unlike traditional adversarial attacks, our loss functions are computed only at a randomly-chosen timestep in the denoising sequence during training. Still, this scheme is effective in breaking the generation output (Sec. 5).

## 4.2. Algorithms

The problem in Eq. 9 is still a challenging bi-level optimization. We define different methods to approximate its solution based on prominent techniques used in literature.

**Fully-trained Surrogate Model Guidance (FSMG).** Most previous image cloaking approaches [49, 49, 57] employ a model trained on clean data as a surrogate to guide the adversarial attack. We can naively follow that direction by using a surrogate DreamBooth model with hyperparameters $\theta_{clean}$ fully finetuned from a small subset of samples $\mathcal{X}_A \subset \mathcal{X}$. This set does not need to cover the target images $\mathcal{X}_{db} = \{x^{(i)}\}_{i=1}^{N_{db}}$; it can be fixed, allowing the surrogate model to be learned once regardless of the constant change of $\mathcal{X}$ and $\mathcal{X}_{db}$. After getting $\theta_{clean}$, we can use it as guidance to find optimal noise for each target image $\delta^{*(i)} = \arg\max_{\delta^{(i)}} \mathcal{L}_{cond}(\theta_{clean}, x^{(i)} + \delta^{(i)})$. By doing so, we expect any DreamBooth model finetuned from the perturbed samples to stay away from $\theta_{clean}$.

**Alternating Surrogate and Perturbation Learning (ASPL).** Using a surrogate model full-trained on clean data may not be a good approximation to solve the problem in Eq. 9. Inspired by [25], we propose to incorporate the training of the surrogate DreamBooth model with the perturbation learning in an alternating manner. The surrogate model $\epsilon_\theta$ is first initiated with the pretrained weights. In each iteration, a clone version $\epsilon'_{\theta'}$ is finetuned on the reference clean data $\mathcal{X}_A$, following Eq. (7). This model is then utilized to expedite the learning of adversarial noises $\delta^{(i)}$ in the current loop. Finally, we update the actual surrogate model $\epsilon_\theta$ on the updated adversarial samples, and move to the next training iteration. We provide a snippet for one training iteration in Eq. 10. With such a procedure, the surrogate model better mimics the true models trained by the malicious DreamBooth users since it is only trained on perturbed data.

$$\theta' \leftarrow \theta.\text{clone}()$$
$$\theta' \leftarrow \arg\min_{\theta'} \sum_{x \in \mathcal{X}_A} \mathcal{L}_{db}(\theta', x)$$
$$\delta^{(i)} \leftarrow \arg\max_{\delta^{(i)}} \mathcal{L}_{cond}(\theta', x^{(i)} + \delta^{(i)}) \qquad (10)$$
$$\theta \leftarrow \arg\min_\theta \sum_{i=1}^{N_{db}} \mathcal{L}_{db}(\theta, x^{(i)} + \delta^{(i)}).$$

**Targeted approaches.** The proposed algorithms above are untargeted; each perturbation noise is learned to maximize the reconstruction loss $\mathcal{L}_{cond}$. Therefore, the adversarial examples $x^{(i)} + \delta^{(i)}$ may guide the target DreamBooth model to learn different adversarial directions, potentially canceling out their effects. Inspired by the success of targeted attacks in [49], we can select a single target $x^{tar}$ to learn optimal $\delta$ such that the output of model is pulled closer to $x_t^{tar}$ when trained with $(x+\delta)_t$. This targeted attack scheme can be plugged into all previous methods, and we denote new algorithms with the prefix "T-", e.g., T-FSMG and T-ASPL.

**Ensemble approaches.** In adverse settings, the pretrained text-to-image generator used by the attacker is unknown. While we can pick one to train the perturbation and hope it

**VGGFace2**

| Method | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|---|---|---|---|---|---|---|---|---|
| | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| No Defense | 0.07 | 0.63 | 0.73 | 15.61 | 0.21 | 0.48 | 0.71 | 9.64 |
| FSMG | 0.56 | **0.33** | **0.31** | **36.61** | 0.62 | 0.29 | 0.37 | 38.22 |
| ASPL | **0.63** | **0.33** | **0.31** | 36.42 | **0.76** | **0.28** | **0.30** | **39.00** |
| T-FSMG | 0.07 | 0.58 | 0.74 | 15.49 | 0.28 | 0.44 | 0.71 | 17.29 |
| T-ASPL | 0.07 | 0.57 | 0.72 | 15.36 | 0.39 | 0.44 | 0.70 | 20.06 |

**CelebA-HQ**

| Method | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|---|---|---|---|---|---|---|---|---|
| | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| No Defense | 0.10 | 0.68 | 0.72 | 17.06 | 0.26 | 0.44 | 0.72 | 7.30 |
| FSMG | **0.34** | **0.48** | 0.56 | 36.13 | **0.35** | **0.36** | 0.66 | 33.60 |
| ASPL | 0.31 | 0.50 | **0.55** | **38.57** | 0.34 | 0.39 | **0.63** | **34.89** |
| T-FSMG | 0.06 | 0.64 | 0.73 | 25.75 | 0.24 | 0.45 | 0.73 | 8.04 |
| T-ASPL | 0.06 | 0.64 | 0.73 | 20.58 | 0.26 | 0.46 | 0.72 | 5.36 |

Table 1: Comparing the defense performance of the proposed methods in a convenient setting.



(a) Comparison between proposed methods
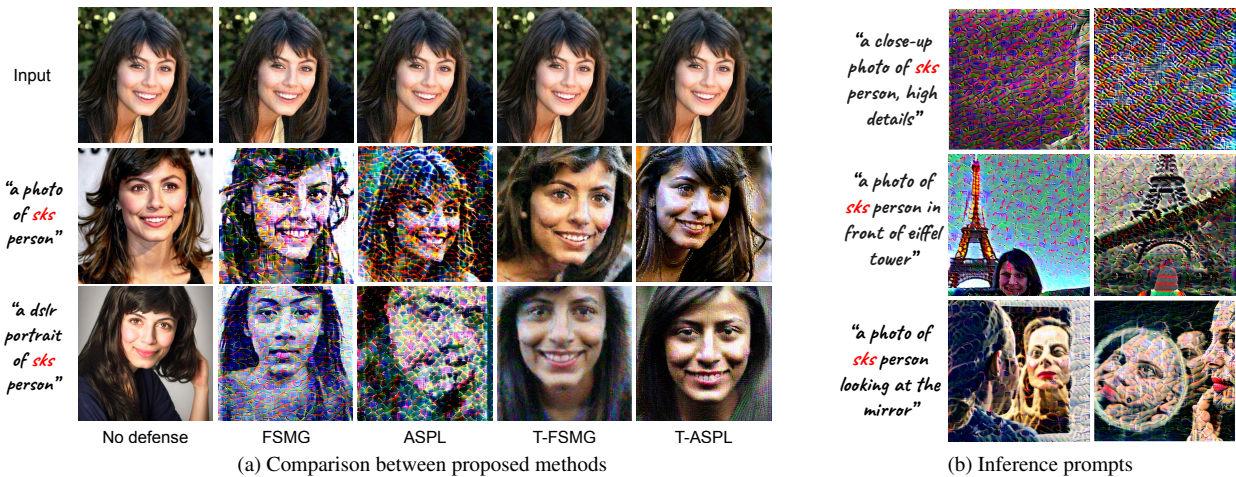
(b) Inference prompts

Figure 2: Qualitative defense results for two subjects in VGGFace2 in the convenient setting. Best viewed in zoom.

transfers well to the target generator, one better approach is to use an ensemble [15, 56] of surrogate models finetuned from different pretrained generators. This approach can be an easy plug-in for the previous approaches. Due to memory constraints, instead of using these surrogate models all at once, we only used a single model at a time, in an interleaving manner, to produce optimal perturbed data.

## 5. Experiments

### 5.1. Experimental setup

**Datasets.** To evaluate the effectiveness of the proposed methods, we look for facial benchmarks that satisfy the following criteria: (1) each dataset covers a large number of different subjects with identity annotated, (2) each subject must have enough images to form two image sets for reference ($\mathcal{X}_A$) and protection ($\mathcal{X}_{db}$), (3) the images should have mid- to high-resolution, (4) the images should be diverse

and in-the-wild. Based on those criteria, we select two famous face datasets CelebA-HQ [28] and VGGFace2 [11].

CelebA-HQ is a high-quality version of CelebA [32] that consists of $30,000$ images at $1024 \times 1024$ resolution. We use the annotated subset [2] that filters and groups images into 307 subjects with at least 15 images for each subject.

VGGFace2 [11] contains around 3.31 million images of 9131 person identities. We filter the dataset to pick subjects that have at least 15 images of resolution above $500 \times 500$.

For fast but comprehensive evaluations, we choose 50 identities for each dataset. For each subject in these datasets, we use the first 12 images and divide them into three subsets, including the reference clean image set, the target protecting set, and an extra clean image set for uncontrolled setting experiments (Sec. 5.5). Each mentioned subset has 4 images with diverse conditions. We then center-crop and resize images to resolution $512 \times 512$.

**Training configurations.** We train each DreamBooth

6

| Version | Defense? | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|---------|----------|------|------|---------|-----------|------|------|---------|-----------|
| | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| v1.4 | ✗ | 0.05 | 0.46 | 0.65 | 21.06 | 0.08 | 0.43 | 0.64 | 10.05 |
| | ✓ | **0.80** | **0.18** | **0.12** | **26.76** | **0.17** | **0.28** | **0.55** | **13.07** |
| v1.5 | ✗ | 0.07 | 0.49 | 0.65 | 18.53 | 0.07 | 0.45 | 0.64 | 10.57 |
| | ✓ | **0.71** | **0.20** | **0.20** | **22.98** | **0.11** | **0.26** | **0.57** | **16.10** |

Table 2: Defense performance of ASPL with different generator versions on VGGFace2 in a convenient setting.

| $\eta$ | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|--------|------|------|---------|-----------|------|------|---------|-----------|
| | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| 0 | 0.07 | 0.63 | 0.73 | 15.61 | 0.21 | 0.48 | 0.71 | 9.64 |
| 0.01 | 0.08 | 0.58 | 0.72 | 33.03 | 0.28 | 0.45 | 0.72 | 17.14 |
| 0.03 | 0.44 | 0.38 | 0.38 | 36.45 | 0.55 | 0.32 | 0.43 | 37.86 |
| 0.05* | 0.63 | 0.33 | 0.31 | 36.42 | 0.76 | 0.28 | 0.30 | 39.00 |
| 0.10 | 0.76 | 0.21 | 0.22 | 37.33 | 0.86 | 0.23 | 0.26 | 40.92 |
| 0.15 | **0.80** | **0.15** | **0.15** | **37.07** | **0.91** | **0.17** | **0.14** | **41.18** |

Table 3: Defense performance of ASPL with different noise budgets on VGGFace2 in a convenient setting. "*" is default.

model, both text-encoder and UNet model, with batch size of 2 and learning rate of $5 \times 10^{-7}$ for 1000 training steps. By default, we use the latest Stable Diffusion (v2.1) as the pretrained generator. Unless specified otherwise, the training instance prompt and prior prompt are "a photo of *sks* person" and "a photo of person", respectively. It takes 15 minutes to train a model on an NVIDIA A100 GPU 40GB.

We optimize the adversarial noise $\delta^{(i)}$ in each step of FSMG and ASPL using the untargeted PGD scheme (Eq. 3). We use 100 PGD iterations for FSMG and 50 iterations for ASPL. Both methods use $\alpha = 0.005$ and the default noise budget $\eta = 0.05$. It takes 2 and 5 minutes for FSMG and ASPL to complete on an NVIDIA A100 GPU 40GB.

**Evaluation metrics.** Our methods aim to disrupt the target DreamBooth models, making them produce poor images of the target user. To measure the defense's effectiveness, for each trained DreamBooth model and each testing prompt, we generate 30 images. We then use a series of metrics to evaluate these generated images comprehensively.

Images generated from successfully disrupted Dream-Booth models may have no detectable face, and we measure that rate, called ***Face Detection Failure Rate (FDFR)***, using RetinaFace detector [17]. If a face is detected, we extract its face recognition embedding, using ArcFace recognizer [18], and compute its cosine distance to the average face embedding of the entire user's clean image set. This metric is called ***Identity Score Matching (ISM)***. Finally, we use two extra image quality assessment metrics. One is ***SER-FQA*** [52], which is an advanced, recent metric dedicated to facial images. The other is ***BRISQUE*** [34], which is classical and popular for assessing images in general.

## 5.2. Convenient setting

We first evaluate the proposed defense methods, including FSMG, ASPL, T-FSMG, and T-ASPL, in a convenient

setting on the two datasets. We try two image generation prompts, one used in training ("a photo of *sks* person") and one novel, unseen prompt ("a dslr portrait of *sks* person"). The average scores over DreamBooth-generated images with each defense are reported in Table 1. As can be seen, the untargeted defenses significantly increase the face detection failure rates and decrease the identity matching scores, implying their success in countering the Dream-Booth threat. We provide some qualitative images in Fig. 2a. As expected, ASPL defends better than FSMG since it mimics better the DreamBooth model training at test time. Targeted methods perform poorly, suggesting that the noise generated by these methods, while providing more consistent adversarial guidance in DreamBooth training, is suboptimal and ineffective. Since ASPL performs the best, we will try only this method in all follow-up experiments.

## 5.3. Ablation studies

**Text-to-image generator version.** In previous experiments, we used Stable Diffusion (SD) v2.1 as the pretrained text-to-image generator. In this section, we examine if our proposed defense (ASPL) is still effective when using different pretrained generators. Since SD is the only open-source large text-to-image model series, we try two of its versions, including v1.4 and v1.5. Note that while belonging to the same model family, these models, including v2.1, are slightly different in networks and behaviors. As reported in Table 2, ASPL shows consistent defense effectiveness.

**Noise budget.** Next, we examine the impact of noise budget $\eta$ on ASPL attack using SD v2.1. As illustrated in Table 3, our defense is already effective with $\eta = 0.03$. The larger the noise budget is, the better defense performance we get, at the cost of the perturbation's stealthiness.

**Inference text prompt.** As indicated in Table 1, ASPL well-disturbs images generated with an unseen text prompt

| | Train | Test | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| Model | v1.4 | v2.1 | 0.62 | 0.31 | **0.28** | 36.00 | 0.70 | 0.31 | 0.35 | 38.39 |
| mismatch | v1.4, 1.5, 2.1 | v2.1 | **0.70** | **0.27** | **0.28** | **36.71** | **0.75** | **0.29** | **0.33** | **39.23** |
| | v1.4 | v2.0 | 0.70 | 0.27 | 0.23 | 36.83 | 0.61 | 0.26 | 0.31 | 37.28 |
| | v1.4, 1.5, 2.1 | v2.0 | **0.79** | **0.24** | **0.18** | **37.96** | **0.71** | **0.23** | **0.23** | **38.99** |
| Term/ Prompt mismatch | DreamBooth prompt | | "a photo of $S_*$ person" | | | | "a dslr portrait of $S_*$ person" | | | |
| | | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| | "*sks*" → "*t@t*" | 0.34 | 0.30 | 0.48 | 36.67 | 0.34 | 0.28 | 0.52 | 28.17 |
| | "a dslr portrait of *sks* person" | 0.07 | 0.15 | 0.69 | 17.34 | 0.49 | 0.37 | 0.36 | 38.42 |

Table 4: Defense performance of ASPL on VGGFace2 when the model, term, or prompt used to train the target DreamBooth model is different from the one used to generate defense noise. Here, $S_*$ is "t@t" for the first row and "sks" for second row.

| Perturbed | Clean | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| 4 | 0 | **0.63** | **0.33** | **0.31** | **36.42** | **0.76** | **0.28** | **0.30** | **39.00** |
| 3 | 1 | 0.50 | 0.43 | 0.41 | 35.53 | 0.52 | 0.35 | 0.51 | 34.01 |
| 2 | 2 | 0.29 | 0.53 | 0.61 | 28.99 | 0.40 | 0.37 | 0.62 | 26.13 |
| 1 | 3 | 0.08 | 0.61 | 0.73 | 18.92 | 0.27 | 0.45 | 0.70 | 15.55 |
| 0 | 4 | 0.07 | 0.63 | 0.73 | 15.61 | 0.21 | 0.48 | 0.71 | 9.64 |

Table 5: Defense performance of ASPL on VGGFace2 in uncontrolled settings. We include two extra results with 0 clean image (convenient setting) and 0 perturbed image (no defense) for comparison.

("a dslr portrait of *sks* person"). We further test the ASPL-disturbed DreamBooth models with different inference text prompts and get similar results. Fig. 2b provides some examples generated with these extra prompts.

### 5.4. Adverse settings

In this section, we investigate if our proposed defense still succeeds when some component is unknown, leading to a mismatch between the perturbation learning and the target DreamBooth model finetuning processes.

**Model mismatching.** The first scenario is when the pre-trained generators are mismatched. We provide an example of transferring adversarial noise trained on SD v1.4 to defend DreamBooth models trained from v2.1 and v2.0 in the first and third rows in Table 4. ASPL still provides good scores as in Table 1. We also examine the ensemble solution suggested in the literature, as discussed in Sec. 4.2. We combine that ensemble idea with ASPL, called E-ASPL, using SD v1.4, 1.5, and 2.1. It further improves the defense in both cases, as illustrated in the upper half of Table 4.

**Term mismatching.** The malicious user can change the term representing the target from the default value ("sks") to another, e.g., "t@t". As reported in the first row in the lower half of Table 4, this term mismatch has only a moderate effect on our results; key scores, like ISM, are still good.

**Prompt mismatching.** The malicious user can also use a different DreamBooth training prompt. ASPL still provides low ISM scores under that adverse scenario, as reported in the last row of Table 4.

**Real-world test.** Our method successfully disrupts person-alized images generated by Astria [1], a black-box commercial service (see the Appendix B for details).

### 5.5. Uncontrolled settings

Our Anti-DreamBooth system is designed for controlled settings, in which all images have protection noises added. In this section, we examine the scenarios when the assumption does not hold, i.e., the malicious user may get in hand some clean images of the target subject and mix them with the perturbed images for DreamBooth training. Assuming the number of images for DreamBooth finetuning is fixed as 4, we examine three data mixing configurations with the number of clean images increasing from 1 to 3. As shown in Table 5, our defense is still quite effective when half of the images are perturbed, but its effectiveness reduces when more clean images are introduced. Still, these uncontrolled settings can be prevented if our system becomes popular and used by all social media with lawmakers' support.

### 6. Conclusions

This paper reveals a potential threat of misused Dream-Booth models and proposes a framework to counter the threat. Our solution is to perturb users' images with subtle adversarial noise so that any DreamBooth model trained on those images will produce poor personalized images. The key idea is to mislead the target DreamBooth model to perform poorly in each denoising step on the original, unperturbed images. We designed several algorithms and extensively evaluate them in different settings. Our defense is ef-

fective, even in adverse conditions. In the future, we aim to improve the perturbation's imperceptibility and robustness [15, 56, 55] and conquer the uncontrolled settings.

# References

[1] Astria. https://www.astria.ai/. 8, 11

[2] CelebA-HQ-Face-Identity-and-Attributes-Recognition-PyTorch. https://github.com/ndb796/CelebA-HQ-Face-Identity-and-Attributes-Recognition-PyTorch. 6

[3] Midjourney. https://www.midjourney.com. 1

[4] Stable Diffusion. https://github.com/Stability-AI/stablediffusion. 2

[5] Fernando Amat, Ashok Chandrashekar, Tony Jebara, and Justin Basilico. Artwork personalization at netflix. In *Proceedings of the 12th ACM Conference on Recommender Systems*, RecSys '18, page 487–488, New York, NY, USA, 2018. Association for Computing Machinery. 2

[6] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search, 2019. 3

[7] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples, 2018. 3

[8] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples, 2017. 3

[9] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 1, 2, 4

[10] Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models, 2017. 3

[11] Qiong Cao, Li Shen, Weidi Xie, Omkar M. Parkhi, and Andrew Zisserman. VGGFace2: A dataset for recognising faces across pose and age. In *International Conference on Automatic Face and Gesture Recognition*, 2018. 6

[12] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57, Los Alamitos, CA, USA, may 2017. IEEE Computer Society. 3

[13] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. 2

[14] Jianbo Chen, Michael I. Jordan, and Martin J. Wainwright. Hopskipjumpattack: A query-efficient decision-based attack, 2019. 3

[15] Valeriia Cherepanova, Micah Goldblum, Harrison Foley, Shiyuan Duan, John P Dickerson, Gavin Taylor, and Tom Goldstein. Lowkey: Leveraging adversarial attacks to protect social media users from facial recognition. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. 3, 6, 9

[16] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks, 2020. 3

[17] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. Retinaface: Single-shot multi-level face localisation in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5203–5212, 2020. 7

[18] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4690–4699, 2019. 7

[19] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 1

[20] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *The Eleventh International Conference on Learning Representations*, 2023. 1, 3

[21] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *CoRR*, abs/1412.6572, 2014. 3

[22] Seok-Ju Hahn, Minwoo Jeong, and Junghye Lee. Connecting low-loss subspace for personalized federated learning. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '22, page 505–515, New York, NY, USA, 2022. Association for Computing Machinery. 2

[23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in neural information processing systems*, 2020. 1, 4

[24] Shengshan Hu, Xiaogeng Liu, Yechao Zhang, Minghui Li, Leo Yu Zhang, Hai Jin, and Libing Wu. Protecting facial privacy: generating adversarial identity masks via style-robust makeup transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15014–15023, 2022. 3

[25] Qidong Huang, Jie Zhang, Wenbo Zhou, Weiming Zhang, and Nenghai Yu. Initiative defense against facial manipulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1619–1627, 2021. 2, 3, 5

[26] David Ingram, Justine Goode, and Anjali Nair. You against the machine: Can you spot which image was created by A.I.? *NBC News*, Dec. 2022. 1

[27] Felix Juefei-Xu, Run Wang, Yihao Huang, Qing Guo, Lei Ma, and Yang Liu. Countering malicious deepfakes: Survey, battleground, and horizon. *International Journal of Computer Vision*, 130(7):1678–1734, 2022. 2

[28] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 6

[29] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. *CVPR*, 2023. 3

[30] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world, 2016. 3

[31] Zheng Li, Ning Yu, Ahmed Salem, Michael Backes, Mario Fritz, and Yang Zhang. Unganable: Defending against gan-based face manipulation. In *USENIX Security*, 2023. 3

[32] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015. 6

[33] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. 3

[34] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. No-reference image quality assessment in the spatial domain. *IEEE Transactions on Image Processing*, 21(12):4695–4708, 2012. 7

[35] S. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2574–2582, Los Alamitos, CA, USA, jun 2016. IEEE Computer Society. 3

[36] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, 2021. 2

[37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1

[38] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. 1

[39] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1

[40] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *ArXiv*, abs/2204.06125, 2022. 2

[41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 1, 2

[42] Kevin Roose. AI-Generated Art Won a Prize. Artists Aren't Happy. *N.Y. Times*, Sept. 2022. 1

[43] Nataniel Ruiz, Sarah Adel Bargal, and Stan Sclaroff. Disrupting deepfakes: Adversarial attacks against conditional image translation networks and facial manipulation systems. In *Computer Vision–ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 236–251. Springer, 2020. 2

[44] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. 2022. 1, 3

[45] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. 1, 2, 4

[46] Axel Sauer, Tero Karras, Samuli Laine, Andreas Geiger, and Timo Aila. Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis. *arXiv preprint arXiv:2301.09515*, 2023. 2

[47] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models. *ArXiv*, abs/2210.08402, 2022. 2

[48] Aviv Shamsian, Aviv Navon, Ethan Fetaya, and Gal Chechik. Personalized federated learning using hypernetworks. In *International Conference on Machine Learning*, pages 9489–9502. PMLR, 2021. 2

[49] Shawn Shan, Emily Wenger, Jiayun Zhang, Huiying Li, Haitao Zheng, and Ben Y Zhao. Fawkes: Protecting privacy against unauthorized deep learning models. In *Proceedings of the 29th USENIX Security Symposium*, 2020. 3, 5

[50] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, 2015. 4

[51] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. 1

[52] Philipp Terhörst, Jan Niklas Kolf, Naser Damer, Florian Kirchbuchner, and Arjan Kuijper. SER-FIQ: unsupervised estimation of face image quality based on stochastic embedding robustness. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 5650–5659. IEEE. 7

[53] Jonathan Uesato, Brendan O'Donoghue, Aaron van den Oord, and Pushmeet Kohli. Adversarial risk and the dangers of evaluating against weak attacks, 2018. 3

[54] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. https://github.com/huggingface/diffusers, 2022. 4

[55] Run Wang, Ziheng Huang, Zhikai Chen, Li Liu, Jing Chen, and Lina Wang. Anti-forgery: Towards a stealthy and robust deepfake disruption attack via adversarial perceptual-aware perturbations. *arXiv preprint arXiv:2206.00477*, 2022. 2, 3, 9

[56] Chaofei Yang, Leah Ding, Yiran Chen, and Hai Li. Defending against gan-based deepfake attacks via transformation-aware adversarial faces. In *2021 international joint conference on neural networks (IJCNN)*, pages 1–8. IEEE, 2021. 2, 3, 6, 9

[57] Chin-Yuan Yeh, Hsi-Wen Chen, Shang-Lun Tsai, and Sheng-De Wang. Disrupting image-translation-based deepfake algorithms with adversarial attacks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision Workshops*, pages 53–62, 2020. 2, 3, 5

[58] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 2

[59] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 1

[60] Yaoyao Zhong and Weihong Deng. Opom: Customized invisible cloak towards face privacy protection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 3

## A. Additional quantitative results

In the main paper, we comprehensively analyzed ASPL's performance on the VGGFace2 dataset. Here, we provide additional quantitative results on the CelebA-HQ dataset. We also report extra results with FSMG, the second-best defense algorithm, on the convenient settings.

### A.1. Ablation studies

**Text-to-image generator version.** We investigate the effectiveness of our defense methods across different versions of SD models, including v1.4 and v1.5.

As reported in Tab. 6, ASPL significantly decreases the identity scores (ISM) in CelebA-HQ, confirming its defense's effectiveness. Its scores, however, are not as good as in VGGFaces2. We can explain it by the fact that CelebA-HQ images are more constrained in pose and quality, reducing the diversity of the image set for DreamBooth and making their combined perturbation effect less severe.

As for FSMG, there is a similar pattern in all metrics on both VGGFace2 and CelebA-HQ, as presented in Tab. 7. FSMG provides a slightly weaker defense compared with ASPL, confirming our observation in the main paper.

**Noise budget.** We further examine the impact of noise budget $\eta$ on FSMG and ASPL using SD v2.1 in Tabs. 8 and 9. As expected, increasing the noise budget leads to better defense scores, either with FSMG or ASPL and either in VGGFace2 or CelebA-HQ. Again, ASPL outperforms FSMG on most evaluation scores.

### A.2. Adverse settings

In the main paper, we verified that our best protection method, i.e., ASPL, remained effective in VGGFace2 when some components of the target DreamBooth training were unknown, resulting in a disparity between the perturbation learning and the DreamBooth finetuning. Here we repeat those defense experiments but on the CelebA-HQ dataset to further confirm ASPL's effectiveness.

**Model mismatching.** As can be seen in Tab. 10, the ASPL approach still works effectively on CelebA-HQ in the cross-model settings. Furthermore, the ensemble approach demonstrates a superior performance on all measurements, the same as the observation on VGGFace2.

**Term mismatching.** In realistic scenarios, the term representing the target in training DreamBooth might vary differently. To demonstrate this problem, we report ASPL's performance when the term "sks" is changed to "t@t". As can be seen in Tab. 10, our method still provides an extremely low ISM score, guaranteeing user protection regardless of the term mismatching.

**Prompt mismatching.** This is the challenging setting when the attacker uses a prompt different from the one used in perturbation learning to train his/her DreamBooth model. In Tab. 10, though there is a drop in some metrics compared with the convenience settings, either the ISM or BRISQUE score remains relatively good. This evidence further assures that our approaches are robust to the prompt mismatching problem.

### A.3. Uncontrolled settings

We examine APSL in the uncontrolled settings on CelebaA-HQ (Tab. 11) and observe the same trend as reported on the VGGFace2 dataset.

## B. Real-world test.

In previous tests, we conducted experiments in laboratory mode. In this section, we examine if our proposed defense actually works in real-world scenarios by trying to disrupt personalized generation outputs of a black-box, commercialized AI service. We find Astria [1] satisfies our criteria and decide to use it in this test. Astria uses the basic DreamBooth setup that allows us to upload images of a specific target subject and input a generation prompt to acquire corresponding synthesized images. It also supports different model settings; we pick the recommended setting (SD v1.5 with face detection enabled) and a totally different one (Protogen 3.4 + Prism) for the tests.

We compare the output of Astria when using the original images and the adversarial images defended by our ASPL method with Stable Diffusion version 2.1 and $\eta = 0.05$ in Figs. 3 and 4, using two different subjects and with each model setting, respectively. As can be seen, our method

significantly reduces the quality of the generated images in various complex prompts and on both target models. Even though these services often rely on proprietary algorithms and architectures that are not transparent to the public, our method remains effective against them. This highlights the robustness of our approach, which can defend against these services without requiring knowledge of their underlying configurations.

## C. Qualitative results

We comprehensively analyzed our defense mechanism quantitatively in the main paper. Here, we provide additional qualitative results to back up those numbers and for visualization, as well.

### C.1. Ablation studies

**Text-to-image generator version.** We compare the defense performance of ASPL using two different versions of SD models (v1.4 and v1.5) on VGGFace2 in Fig. 5. The output images produced by both models with both prompts are strongly distorted with notable artifacts. We observe the same behavior in the corresponding experiments on CelebA-HQ, visualized in Fig. 6.

**Noise budget.** In order to better understand the impact of the noise budget, we present a grid of images for ASPL on VGGFace2 where the upper bound of noise's magnitude increases along the vertical axis in Fig. 7. It is evident that when the noise budget increases, the visibility of noise becomes more pronounced. Moreover, the allocated noise budget heavily influences the degree of output distortion, resulting in a trade-off between the visibility of noise in perturbed images and the level of output distortion. For further visulization on CelebA-HQ, please refer to Fig. 8.

### C.2. Adverse Setting

**Model mismatching.** In this section, we present the visual outputs of ASPL when a model mismatch occurs. Specifically, we train the image perturbation with SD v1.4, then use those images to disrupt DreamBooth models finetuned from v2.1 and v2.0, respectively. As illustrated in Fig. 9, our defense method is still effective in both cases, although transferring from v1.4 to v2.0 produces more noticeable artifacts than the previous scenario.

In addition to our primary analysis, our study provides qualitative results for E-ASPL, which employs an ensemble method to overcome the challenge of model mismatching. Specifically, we combined knowledge from three versions of SD models (v1.4, v1.5, and v2.1). The results, illustrated in Fig. 10, demonstrate the superior performance of E-ASPL in countering model mismatching where most images are heavily distorted.

**Term mismatching.** Despite the discrepancy of term replacement (from "sks" to "t@t"), ASPL still demonstrates its effectiveness on two provided subjects and two provided prompts (as in Fig. 11). However, the change in the training term may result in slightly weaker artifacts compared to the original setting.

**Prompt mismatching.** The results depicted in Fig. 11 indicate that the finetuning of the DreamBooth model with various prompts, such as "a DSLR portrait of *sks* person", can impact the degree of output distortion to some extent. It is important to note that prompt mismatching can alter the behavior of our defense method on a different prompt, such as "a photo of *sks* person", which can change the identity of the target subject in the generated images.

### C.3. Uncontrolled settings.

All previous results are for controlled settings, in which we have access to all images needing protection. Here, we also include some qualitative results for uncontrolled settings where a mixture of clean and perturbed images are used for finetuning Dreambooth. We use the same settings as the one in the main paper, with the number of images for DreamBooth being fixed at 4 and the number of clean images gradually increase. As can be seen in Fig. 12, our method is more effective when more perturbed data are used and vice versa.

| Version | Defense? | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| v1.4 | ✗ | 0.07 | 0.48 | 0.66 | 16.09 | **0.11** | 0.40 | 0.67 | 10.31 |
| | ✓ | **0.28** | **0.29** | **0.47** | **20.05** | 0.06 | **0.31** | **0.64** | **10.55** |
| v1.5 | ✗ | 0.06 | 0.53 | 0.69 | 14.45 | **0.07** | 0.39 | 0.68 | 8.95 |
| | ✓ | **0.16** | **0.36** | **0.58** | **21.09** | 0.06 | **0.26** | **0.64** | **12.28** |

Table 6: Defense performance of ASPL with different generator versions on CelebA-HQ in a convenient setting.

| VGGFace2 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Version | Defense? | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
| | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| v1.4 | ✗ | 0.05 | 0.46 | 0.65 | 21.06 | 0.08 | 0.44 | 0.64 | 10.05 |
| | ✓ | **0.73** | **0.21** | **0.17** | **25.88** | **0.13** | **0.28** | **0.57** | **13.46** |
| v1.5 | ✗ | 0.07 | 0.49 | 0.65 | 18.53 | 0.07 | 0.45 | 0.64 | 10.57 |
| | ✓ | **0.61** | **0.21** | **0.26** | **23.89** | **0.11** | **0.26** | **0.57** | **18.00** |
| **CelebA-HQ** | | | | | | | | | |
| Version | Defense? | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
| | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| v1.4 | ✗ | 0.07 | 0.48 | 0.66 | 16.09 | **0.11** | 0.40 | 0.67 | 10.31 |
| | ✓ | **0.29** | **0.32** | **0.48** | **20.83** | 0.07 | **0.29** | **0.63** | **12.00** |
| v1.5 | ✗ | 0.06 | 0.53 | 0.69 | 14.45 | **0.07** | 0.39 | 0.68 | 8.95 |
| | ✓ | **0.13** | **0.38** | **0.60** | **20.43** | 0.06 | **0.28** | **0.65** | **13.27** |

Table 7: Defense performance of FSMG with different generator versions on VGGFace2 and CelebA-HQ in a convenient setting.

| VGGFace2 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $\eta$ | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
| | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| 0 | 0.07 | 0.63 | 0.73 | 15.61 | 0.21 | 0.48 | 0.71 | 9.64 |
| 0.01 | 0.09 | 0.58 | 0.73 | 31.58 | 0.28 | 0.46 | 0.71 | 15.85 |
| 0.03 | 0.45 | 0.39 | 0.38 | **37.82** | 0.53 | 0.33 | 0.47 | 38.17 |
| 0.05* | 0.56 | 0.33 | 0.31 | 36.61 | 0.62 | 0.29 | 0.37 | 38.22 |
| 0.10 | 0.70 | 0.22 | 0.23 | 36.60 | 0.77 | 0.27 | 0.29 | 38.59 |
| 0.15 | **0.77** | **0.20** | **0.20** | 36.16 | **0.83** | **0.22** | **0.26** | **39.17** |
| **CelebA-HQ** | | | | | | | | |
| $\eta$ | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
| | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| 0 | 0.10 | 0.68 | 0.72 | 17.06 | 0.26 | 0.44 | 0.72 | 7.30 |
| 0.01 | 0.12 | 0.68 | 0.73 | 19.55 | 0.30 | 0.46 | 0.71 | 6.60 |
| 0.03 | 0.15 | 0.57 | 0.71 | 33.89 | 0.27 | 0.41 | 0.73 | 22.67 |
| 0.05* | 0.34 | 0.48 | 0.56 | 36.13 | 0.35 | 0.36 | 0.66 | 33.60 |
| 0.10 | 0.73 | 0.32 | 0.27 | **39.16** | 0.67 | 0.24 | 0.43 | **38.99** |
| 0.15 | **0.77** | **0.29** | **0.26** | 38.22 | **0.73** | **0.23** | **0.35** | 38.22 |

Table 8: Defense performance of FSMG with different noise budgets on VGGFace2 and CelebA-HQ in a convenient setting. "*" is default.

| $\eta$ | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|---|---|---|---|---|---|---|---|---|
| | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| 0 | 0.10 | 0.68 | 0.72 | 17.06 | 0.26 | 0.44 | 0.72 | 7.30 |
| 0.01 | 0.11 | 0.67 | 0.72 | 19.97 | 0.27 | 0.45 | 0.72 | 6.65 |
| 0.03 | 0.12 | 0.60 | 0.71 | 34.34 | 0.25 | 0.44 | 0.73 | 18.29 |
| 0.05* | 0.31 | 0.50 | 0.55 | 38.57 | 0.34 | 0.39 | 0.63 | 34.89 |
| 0.10 | 0.73 | 0.36 | 0.30 | **38.83** | 0.74 | 0.27 | 0.36 | **38.96** |
| 0.15 | **0.86** | **0.25** | **0.19** | 38.67 | **0.82** | **0.24** | **0.28** | 38.86 |

Table 9: Defense performance of ASPL with different noise budgets on CelebA-HQ in a convenient setting. "*" is default.

| | Train | Test | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| Model mismatch | v1.4 | v2.1 | 0.37 | 0.48 | 0.53 | 39.28 | 0.34 | 0.39 | 0.64 | 33.50 |
| | v1.4, 1.5, 2.1 | v2.1 | **0.39** | **0.46** | **0.48** | **38.25** | **0.44** | **0.34** | **0.57** | **37.29** |
| | v1.4 | v2.0 | 0.40 | 0.46 | 0.51 | 38.88 | 0.43 | 0.36 | 0.60 | 22.21 |
| | v1.4, 1.5, 2.1 | v2.0 | **0.56** | **0.43** | **0.43** | **41.83** | **0.55** | **0.33** | **0.51** | **29.93** |
| Term/ Prompt mismatch | DreamBooth prompt | | "a photo of $S_*$ person" | | | | "a dslr portrait of $S_*$ person" | | | |
| | | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| | "*sks*" → "*t@t*" | | 0.20 | 0.17 | 0.64 | 26.49 | 0.17 | 0.10 | 0.65 | 1.14 |
| | "a dslr portrait of *sks* person" | | 0.13 | 0.22 | 0.69 | 18.51 | 0.33 | 0.51 | 0.58 | 37.99 |

Table 10: Defense performance of ASPL on CelebA-HQ when the model, term, or prompt used to train the target DreamBooth model is different from the one used to generate defense noise. Here, $S_*$ is "t@t" for the first row and "sks" for second row.

| Perturbed | Clean | "a photo of *sks* person" | | | | "a dslr portrait of *sks* person" | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ | FDFR↑ | ISM↓ | SER-FQA↓ | BRISQUE↑ |
| 4 | 0 | **0.31** | **0.50** | **0.55** | **38.57** | **0.34** | **0.39** | **0.63** | **34.89** |
| 3 | 1 | 0.26 | 0.54 | 0.63 | 32.23 | 0.30 | 0.40 | 0.69 | 22.03 |
| 2 | 2 | 0.19 | 0.61 | 0.69 | 25.14 | 0.25 | 0.41 | 0.71 | 11.35 |
| 1 | 3 | 0.13 | 0.65 | 0.72 | 19.24 | 0.23 | 0.43 | 0.72 | 9.70 |
| 0 | 4 | 0.10 | 0.68 | 0.72 | 17.06 | 0.26 | 0.44 | 0.72 | 7.30 |

Table 11: Defense performance of ASPL on CelebA-HQ in uncontrolled settings. We include two extra results with 0 clean image (convenient setting) and 0 perturbed image (no defense) for comparison.
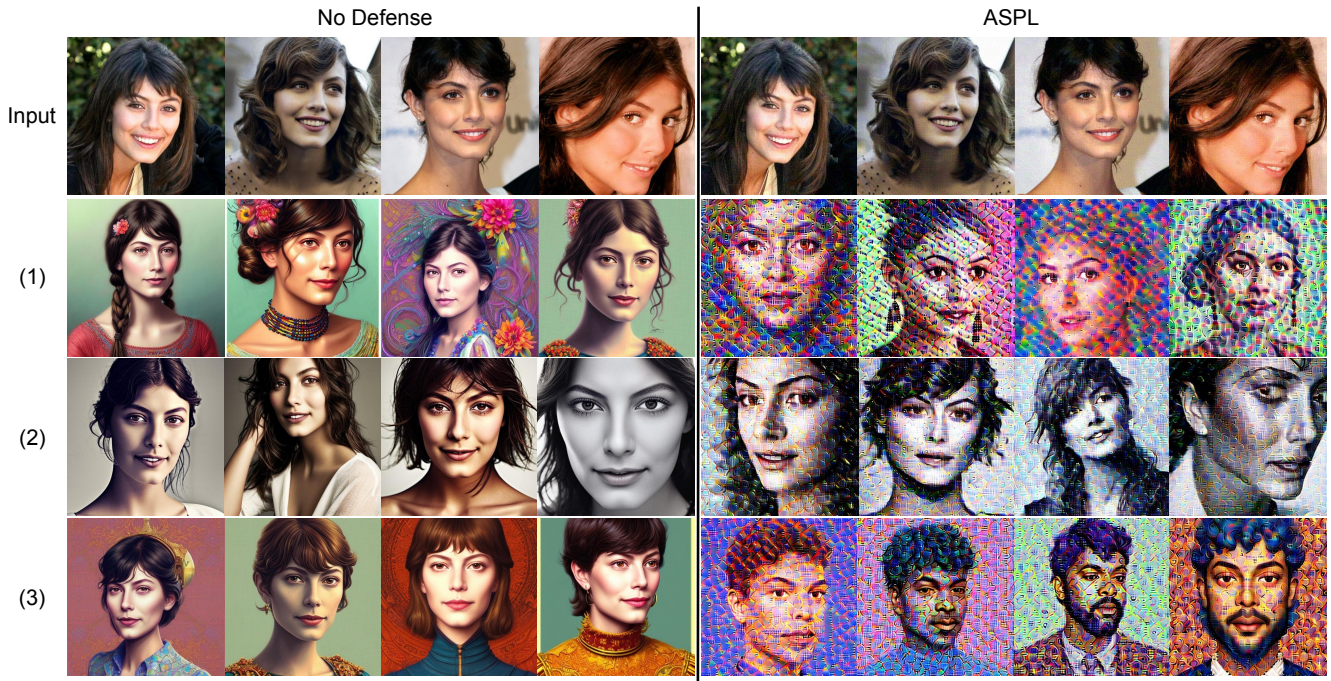
Figure 3: **Disrupting personalized images generated by Astria (SD v1.5 with face detection enabled)**. The prompts for image generation include: (1) "portrait of *sks* person portrait wearing fantastic Hand-dyed cotton clothes, embellished beaded feather decorative fringe knots, colorful pigtail, subtropical flowers and plants, symmetrical face, intricate, elegant, highly detailed, 8k, digital painting, trending on pinterest, harper's bazaar, concept art, sharp focus, illustration, by artgerm, Tom Bagshaw, Lawrence Alma-Tadema, greg rutkowski, alphonse Mucha", (2) "close up of face of *sks* person fashion model in white feather clothes, official balmain editorial, dramatic lighting highly detailed", and (3) "portrait of sks person prince :: by Martine Johanna and Simon Stålenhag and Chie Yoshii and Casey Weldon and wlop :: ornate, dynamic, particulate, rich colors, intricate, elegant, highly detailed, centered, artstation, smooth, sharp focus, octane render, 3d"
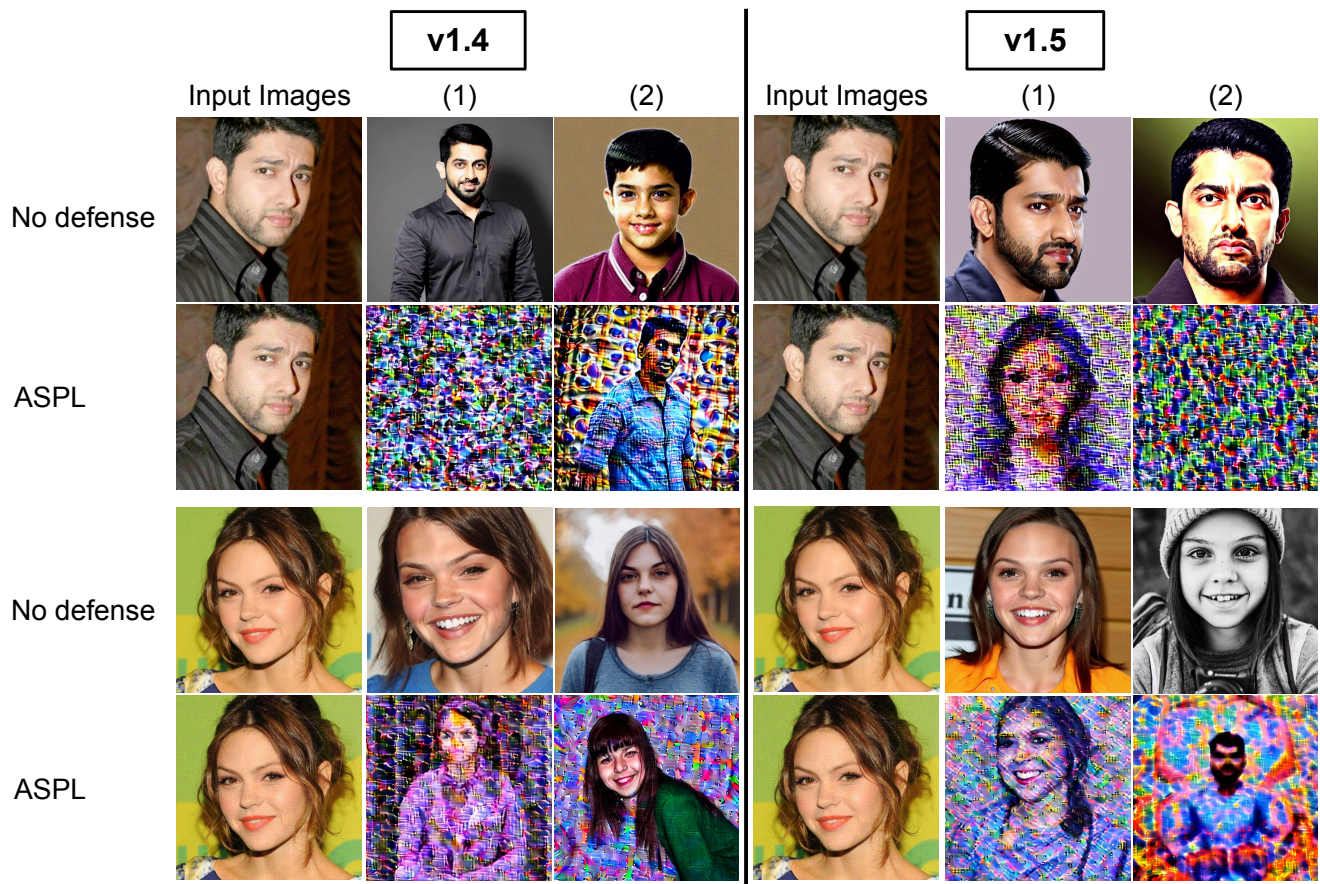
Figure 4: **Disrupting personalized images generated by Astria (Protogen with Prism and face detection enabled)**. The prompts for image generation include: (1) "portrait of *sks* person portrait wearing fantastic Hand-dyed cotton clothes, embellished beaded feather decorative fringe knots, colorful pigtail, subtropical flowers and plants, symmetrical face, intricate, elegant, highly detailed, 8k, digital painting, trending on pinterest, harper's bazaar, concept art, sharp focus, illustration, by artgerm, Tom Bagshaw, Lawrence Alma-Tadema, greg rutkowski, alphonse Mucha", (2) "close up of face of *sks* person fashion model in white feather clothes, official balmain editorial, dramatic lighting highly detailed", and (3) "portrait of sks person prince :: by Martine Johanna and Simon Stålenhag and Chie Yoshii and Casey Weldon and wlop :: ornate, dynamic, particulate, rich colors, intricate, elegant, highly detailed, centered, artstation, smooth, sharp focus, octane render, 3d"

Figure 5: Qualitative results of ASPL with two different versions of SD models (v1.4 and v1.5) on VGGFace2. We provide in each test a single, representative input image. The generation prompts include (1) "a photo of *sks* person" and (2) "a dslr portrait of *sks* person".
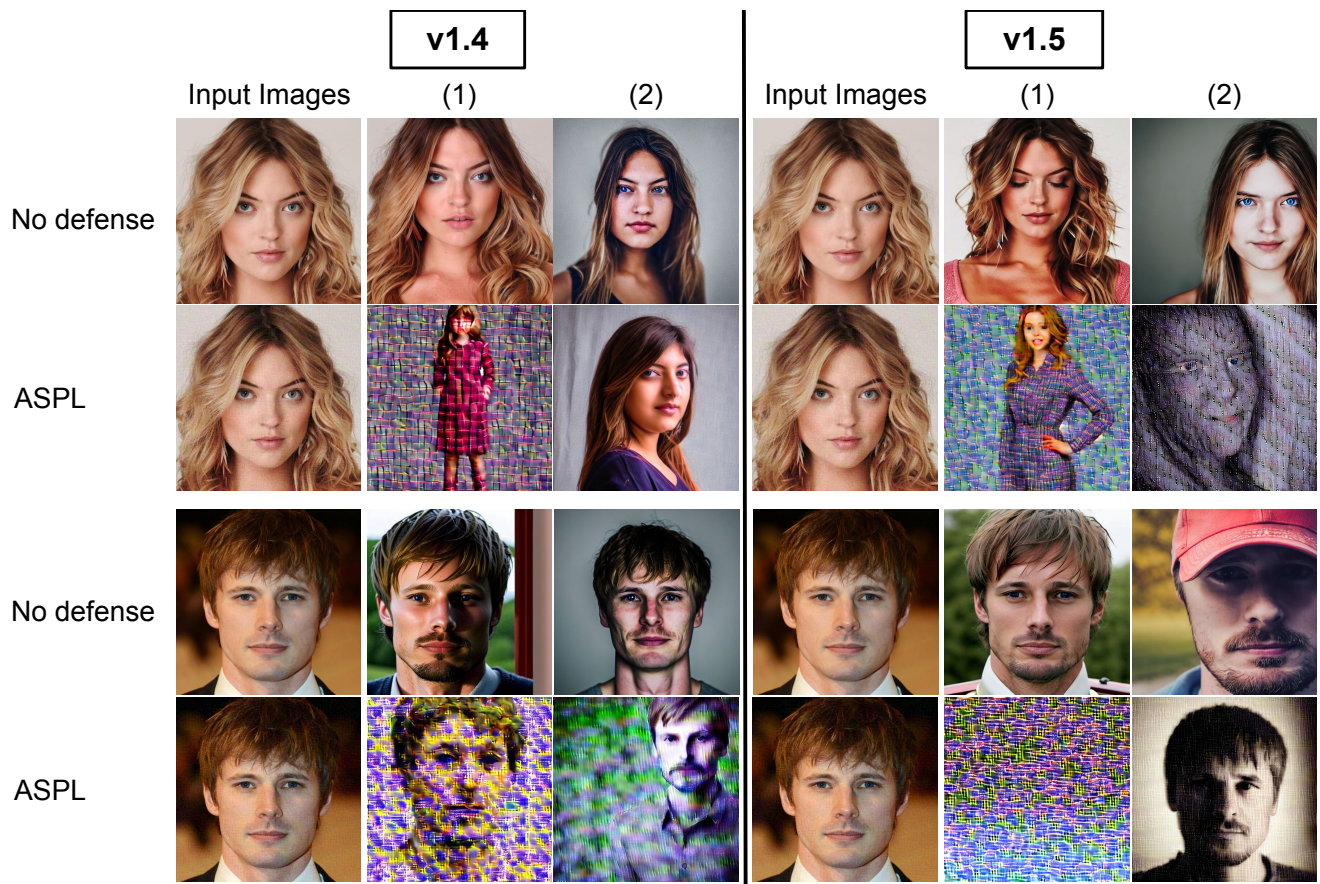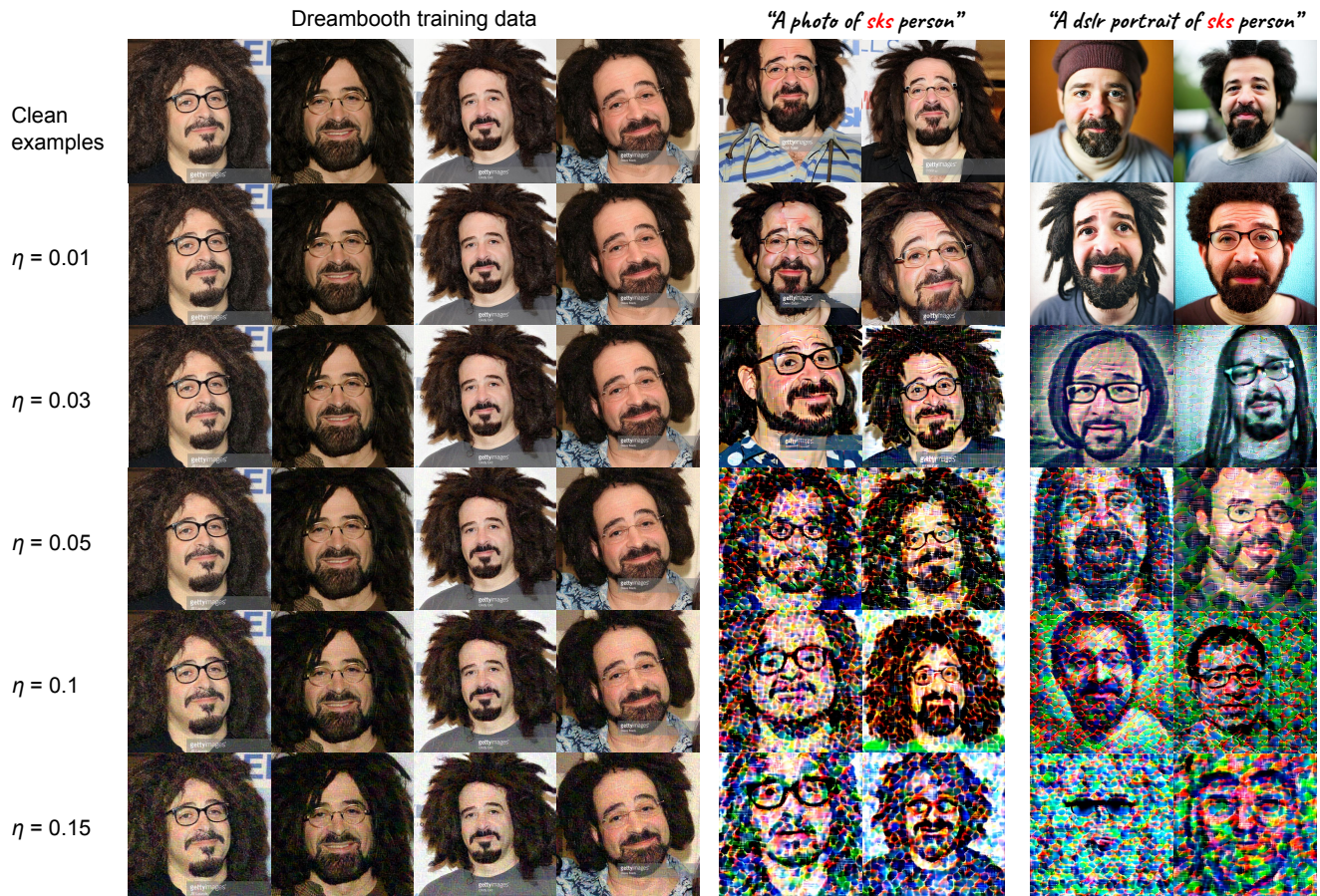
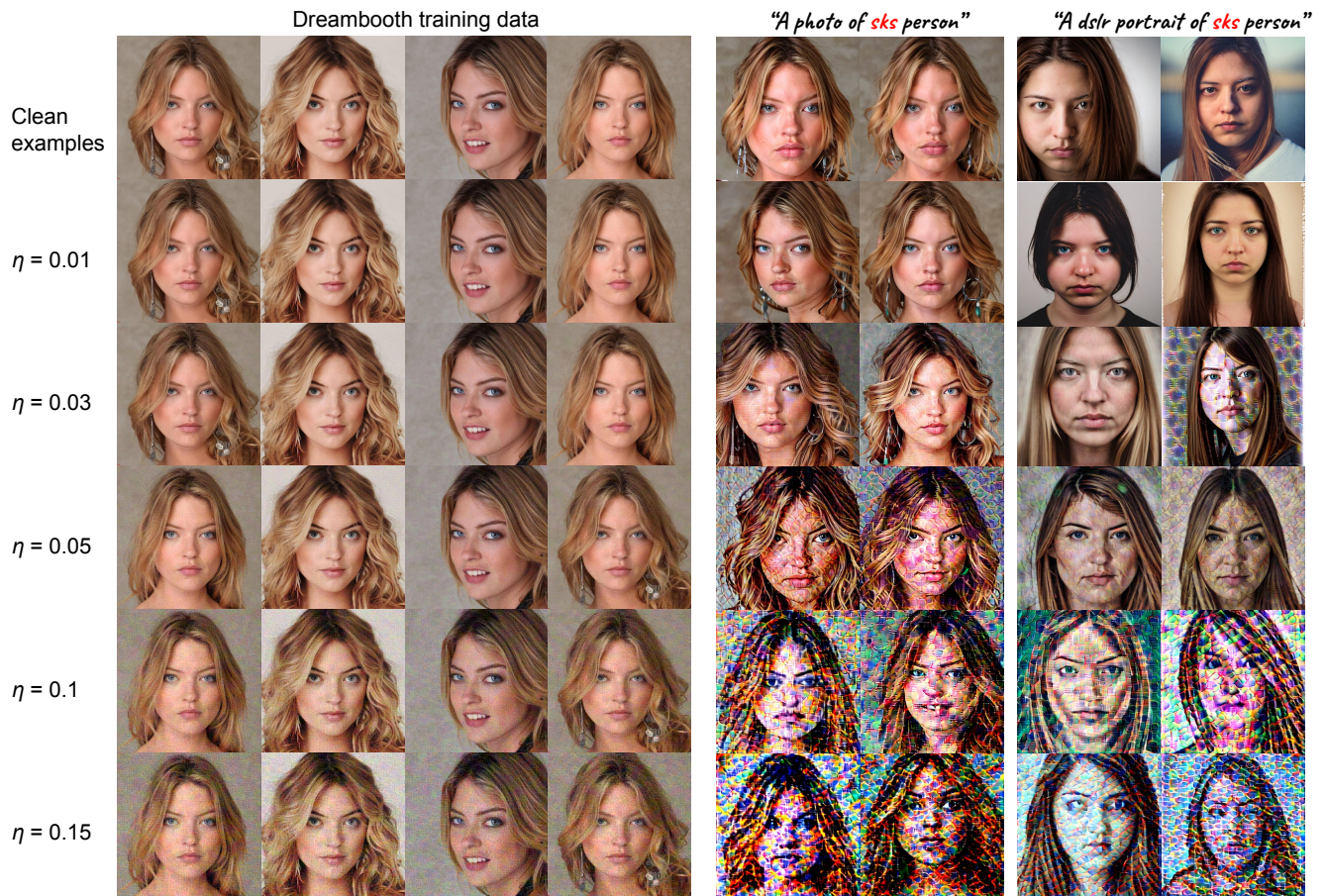Figure 6: Qualitative results of ASPL with two different versions of SD models (v1.4 and v1.5) on CelebA-HQ. We provide in each test a single, representative input image. The generation prompts include (1) "a photo of *sks* person" and (2) "a dslr portrait of *sks* person".

Figure 7: Qualatative results of ASPL with different noise budget on VGGFace2.

Figure 8: Qualitative results of ASPL with different noise budget on CelebA-HQ.

Figure 9: Qualitative results of ASPL in adverse settings on VGGFace2 where the SD model version in perturbation learning mismatches the one used in the DreamBooth finetuning stage (v1.4 → v2.1 and v1.4 → v2.0). We test with two random subjects and denote them in green and red, respectively.
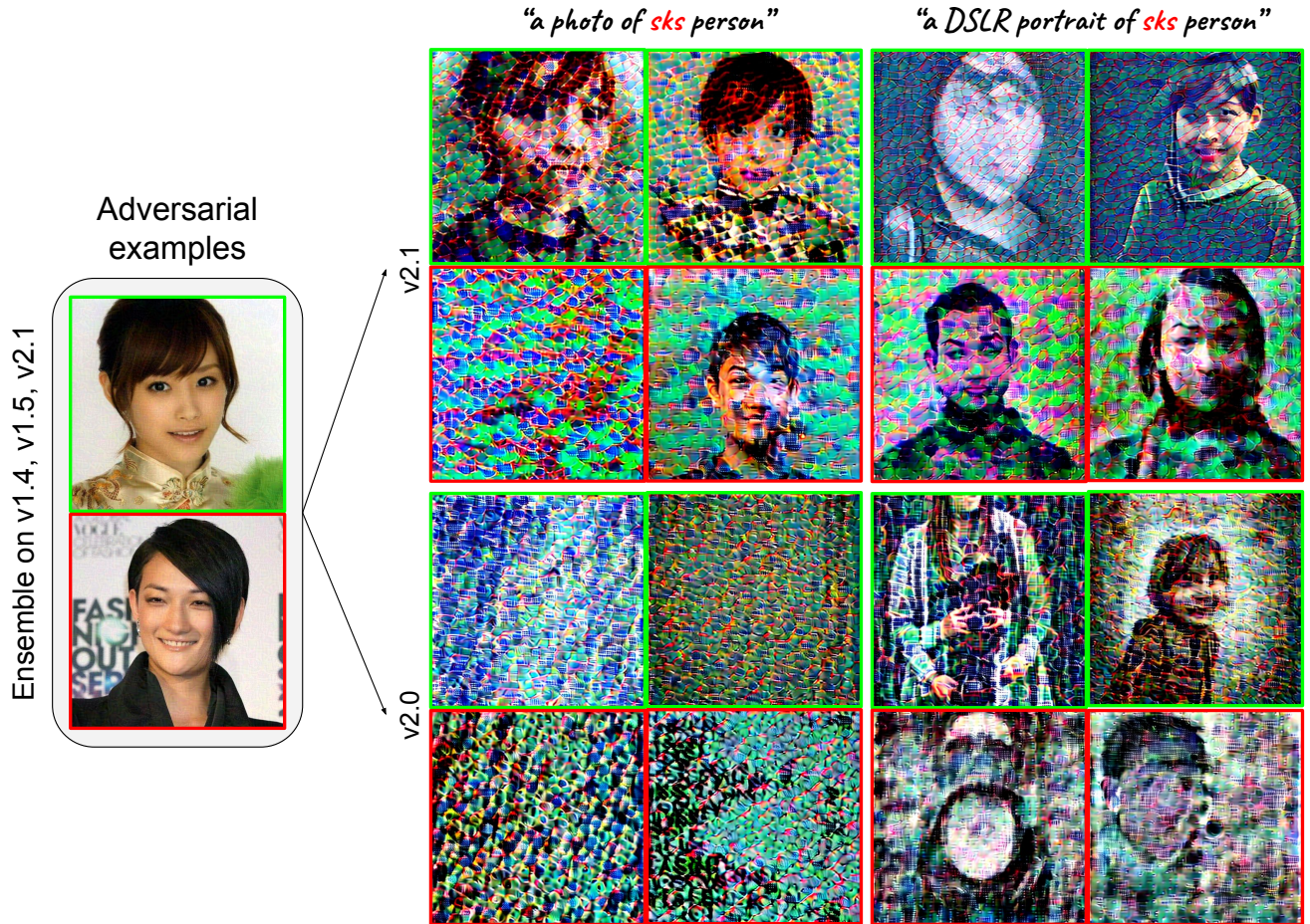
Figure 10: Qualitative results of E-ASPL on VGGFace2, where the ensemble model combines 3 versions of SD models, including v1.4, v1.5, and v2.1. Its performance is validated on two DreamBooth models finetuned on SD v2.1 and v2.0, respectively. We test with two random subjects and denote them in green and red, respectively.
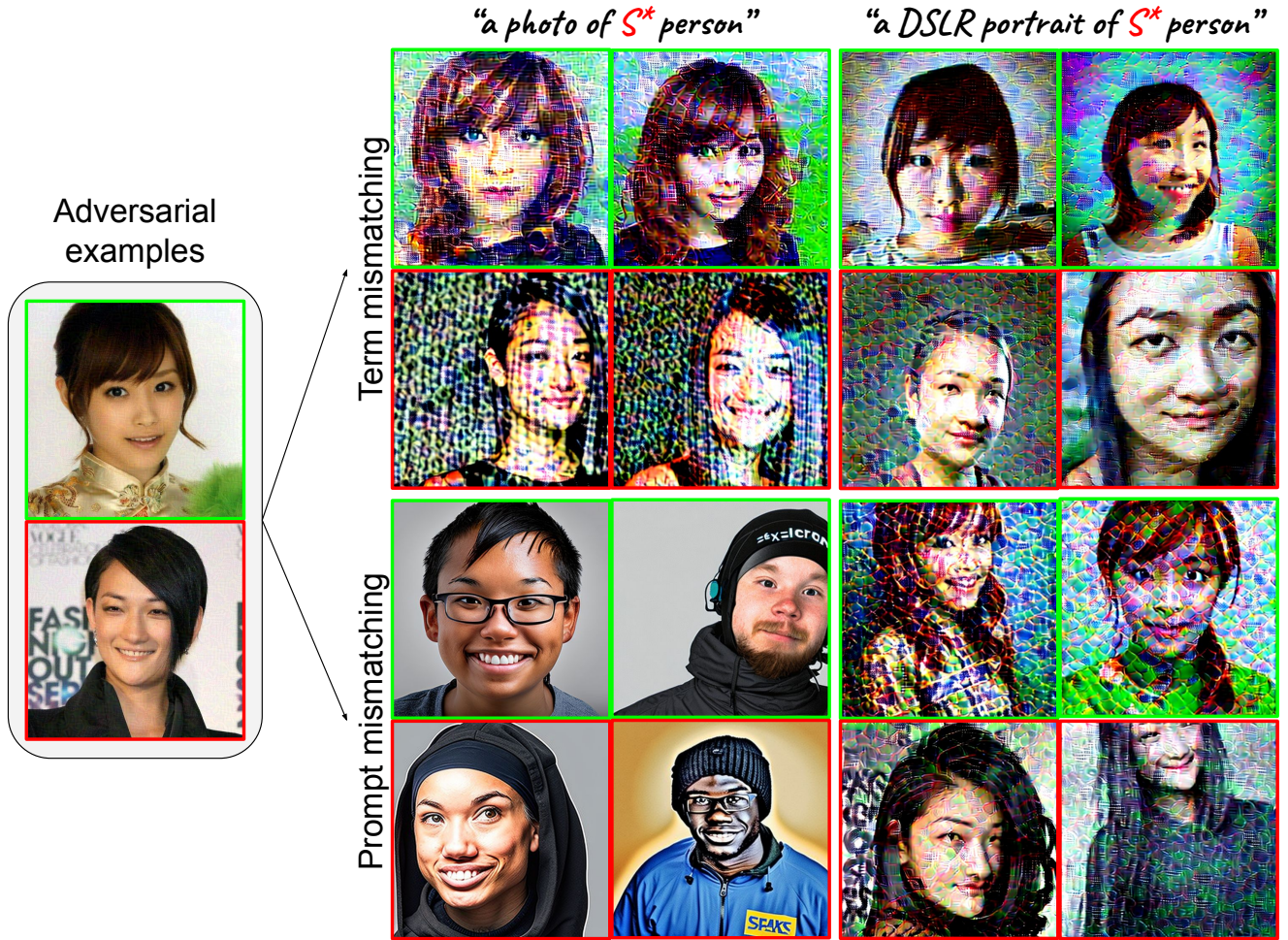
Figure 11: Qualitative results of ASPL on VGGFace2 where the training term and prompt of the target DreamBooth model mismatch the ones in perturbation learning. In the first scenario, the training term is changed from "sks" to "t@t". In the second scenario, the training prompt is replaced with "a DSLR portrait of *sks* person" instead of "a photo of *sks* person". Here, $S^*$ is "t@t" for term mismatching and "sks" for prompt mismatching. We test with two random subjects and denote them in green and red, respectively.
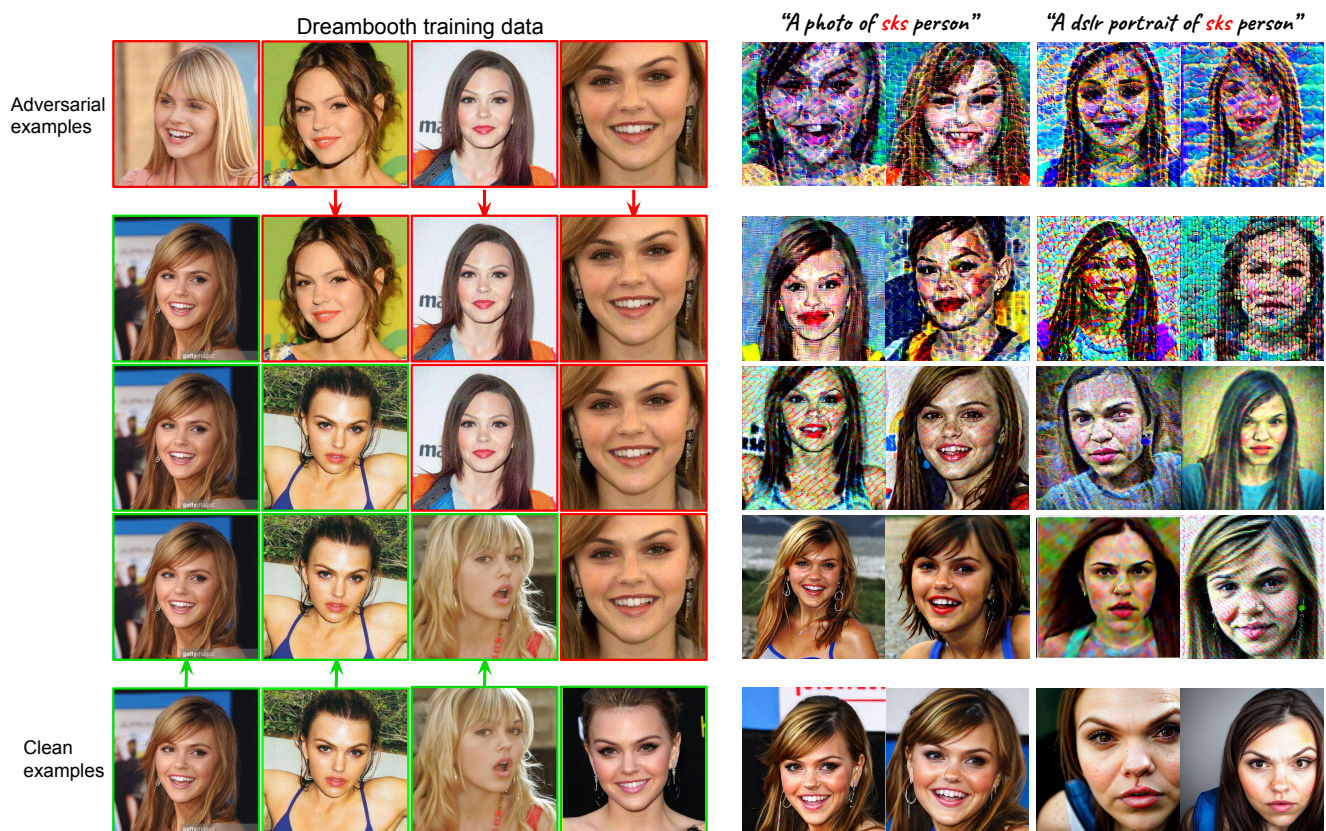
Figure 12: Qualatitive results of ASPL in uncontrolled setting on VGGFace2. We denote the perturbed examples and the leaked clean examples in red and green, respectively.