# FreezeAsGuard: Mitigating Illegal Adaptation of Diffusion Models via Selective Tensor Freezing

**Kai Huang, Haoming Wang and Wei Gao**
University of Pittsburgh
k.huang, hw.wang, weigao@pitt.edu

## Abstract

Text-to-image diffusion models can be fine-tuned in custom domains to adapt to specific user preferences, but such adaptability has also been utilized for illegal purposes, such as forging public figures' portraits, duplicating copyrighted artworks and generating explicit contents. Existing work focused on detecting the illegally generated contents, but cannot prevent or mitigate illegal adaptations of diffusion models. Other schemes of model unlearning and reinitialization, similarly, cannot prevent users from relearning the knowledge of illegal model adaptation with custom data. In this paper, we present *FreezeAsGuard*, a new technique that addresses these limitations and enables irreversible mitigation of illegal adaptations of diffusion models. Our approach is that the model publisher selectively freezes tensors in pre-trained diffusion models that are critical to illegal model adaptations, to mitigate the fine-tuned model's representation power in illegal adaptations, but minimize the impact on other legal adaptations. Experiment results in multiple text-to-image application domains show that FreezeAsGuard provides 37% stronger power in mitigating illegal model adaptations compared to competitive baselines, while incurring less than 5% impact on legal model adaptations. The source code is available at: https://github.com/pittisl/FreezeAsGuard.

## 1 Introduction

Text-to-image diffusion models [44, 43] are powerful tools to generate high-quality images aligned with user prompts. After pre-trained by model publishers to embed world knowledge from large image data [49], open-sourced diffusion models, such as Stable Diffusion (SD) [9, 10], can be conveniently adapted by users to generate their preferred images[1], through fine-tuning with custom data in specific domains. For example, diffusion models can be fine-tuned on cartoon datasets to synthesize avatars in video games [46], or on datasets of landscape photos to generate wallpapers [11].

An increasing risk of democratizing open-sourced diffusion models, however, is that the capability of model adaptation has been utilized for illegal purposes, such as forging public figures' portraits [22, 24], duplicating copyrighted artworks [26], and generating explicit content [25]. Most existing efforts aim to deter at-
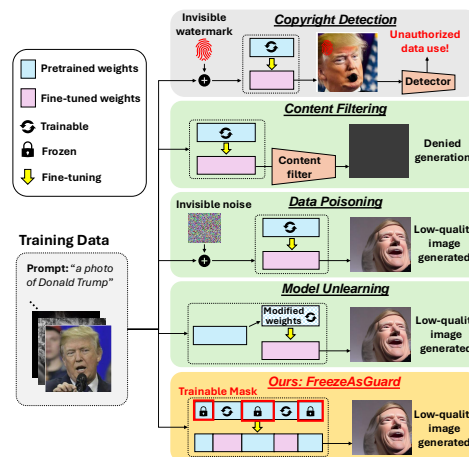


Figure 1: Existing work vs. FreezeAsGuard in mitigating malicious adaptation of diffusion models

---

[1]Many APIs, such as HuggingFace Diffusers [56], can be used for fine-tuning open-sourced diffusion models with the minimum user efforts.

Figure 2: FreezeAsGuard ensures that portraits (left) and artworks (right) generated by diffusion models in illegal classes cannot be recognizable as target objects, even if the model has been fine-tuned with data samples in illegal classes. In contrast, unlearning schemes (UCE [23] and IMMA [65]) cannot prevent the unlearned knowledge of illegal classes from being relearned in fine-tuning.

tempts of illegal model adaptation with copyright detection [64, 16, 17], which embeds invisible but detectable watermarks into training data and further generated images, as shown in Figure 1. However, such detection only applies to misuse of training data, and does not mitigate the user's capability of illegal model adaptation. Users can easily bypass such detection by collecting and using their own training data without being watermarked (e.g., users' self-taken photos of public figures).

Instead, an intuitive approach to mitigation is content filtering. However, filtering user prompts [19] can be bypassed by fine-tuning the model to align innocent prompts with illegal image contents [55], and filtering the generated images [7] is often overpowered with high false-positive rates [3]. Data poisoning techniques can avoid false positives by injecting invisible perturbations into training data [59, 62, 51], but cannot apply when public web data or users' private data is used for fine-tuning. Recent unlearning methods allow model publishers to remove knowledge needed for illegal adaptation by modifying model weights [20, 23, 57, 65] , but cannot prevent relearning such knowledge via fine-tuning.

The key limitation of these techniques is that they focus on modifying the training data or model weights, but such modification can be reversed by users via fine-tuning with their own data. Such modification, further, cannot restrain the mitigation power only in *illegal data classes* (e.g., public figures' portraits) without affecting model adaptation in other *legal data classes* (e.g., the user's own portraits), due to the high ambiguity and possible overlap between these classes.

To prevent users from reversing the mitigation maneuvers being applied, in this paper we present *FreezeAsGuard*, a new technique that constrains the trainability of diffusion model's tensors in fine-tuning. As shown in Figure 1, the model publisher selectively freezes tensors in pre-trained models that are critical to fine-tuning in illegal classes (e.g., public figures' portraits), to limit the model's representation power of being fine-tuned in illegal classes. In practice, since most illegal users are not professional and fine-tune diffusion models by simply following the instructions provided by model publishers, tensor freezing can be effectively enforced by model publishers through these instructions, to guide the users to adopt tensor freezing. Essentially, since freezing tensors lowers the trainable model parameters and reduces the computing costs of fine-tuning, users would be well motivated to adopt tensor freezing in fine-tuning practices.

The major challenge is how to properly evaluate the importance of tensors in model fine-tuning. Popular attribution-based importance metrics [38, 41] are used in model pruning with fixed weight values, but cannot reflect the impact of weight variations in fine-tuning. Such impact of weight variations, in fact, cannot be condensed into a single importance metric, due to the randomness and interdependencies of weight updates in fine-tuning iterations.



Figure 3: Mask learning and fine-tuning as a bilevel optimization

Instead, as shown in Figure 3, we formulate the selection of frozen tensors in all the illegal classes as one *trainable binary mask*. Given a required ratio
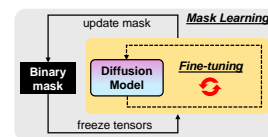
of frozen tensors specified by model publisher, we optimize such selection
with training data in all the involved illegal classes, through bilevel optimization that combines the
iterative process of mask learning and iterations of model fine-tuning. In this way, the mask being
trained can timely learn the impact of weight variations on the training loss during fine-tuning.

With frozen tensors, the model's representation power should be retained when fine-tuned on other
legal classes (e.g., user's own portraits). Hence, we incorporate training samples from legal classes into
the bilevel optimization, to provide suppressing signals for selecting tensors being frozen. Hence, the
learned mask of freezing tensors should skip tensors that are important to fine-tuning in legal classes.

We evaluated FreezeAsGuard in three different domains of illegal model adaptations: *1)* forging public
figures' portraits, *2)* duplicating copyrighted artworks and *3)* generating explicit contents. For each
domain, we use open-sourced or self-collected datasets, and randomly select different data classes as
illegal and legal classes. We use competitive model unlearning schemes as baselines, and multiple
metrics to measure image quality. Our findings are as follows:

- FreezeAsGuard has strong mitigation power in illegal classes. Compared to the competitive
  baselines, it further reduces the quality of images generated by fine-tuned model by up to
  37%, and ensures the generated images to be unrecognizable as subjects in illegal classes.
- FreezeAsGuard has the minimum impact on modal adaptation in legal classes. It ensures
  on-par quality of the generated images compared to regular full fine-tuning on legal data, with
  a difference of at most 5%.
- FreezeAsGuard has high compute efficiency. Compared to full fine-tuning, it can save up to
  48% GPU memory and 21% wall-clock computing time.

## 2  Background & Motivation

### 2.1  Fine-Tuning Diffusion Models

Given text prompts $y$ and images $x$ as training data, fine-tuning a diffusion model approximates the
conditional distribution $p(x|y)$ by learning to reconstruct images that are progressively blurred with
noise $\epsilon$ over step $t = 1, ..., T$. Training objective is to minimize the reconstruction loss:

$$\mathcal{L}_\theta = \mathbb{E}_{x,y,\epsilon \sim \mathcal{N}(0,1),t} \left[ \|\epsilon - \epsilon_\theta(\mathcal{E}(x_t), t, \tau(y))\|_2^2 \right], \tag{1}$$

where $\mathcal{E}(\cdot)$ is the encoder of a pretrained VAE, $\tau(\cdot)$ is a pretrained text encoder, and $\epsilon_\theta(\cdot)$ is a denoising
model with trainable parameters $\theta$. Most diffusion models adopt UNet architecture [45] as the denoising
model.

In fine-tuning, the diffusion model learns new knowledge by adapting the generic knowledge in the
pre-trained model [13]. For example, new knowledge about "a green beetle" can be a combination
of generic knowledge on "hornet" and "emerald". This behavior implies that fine-tuning in different
classes may share the same knowledge base, and it is challenging to focus the mitigation power in
illegal classes without affecting fine-tuning in other legal classes. This challenge motivates us to
regulate FreezeAsGuard's mitigation power by incorporating training samples in legal classes, when
selecting tensors being frozen for illegal classes.

| Model component Being frozen | CLIP ($\uparrow$) | TOPIQ ($\uparrow$) | FID ($\downarrow$) |
|---|---|---|---|
| No freezing | 31.93 | 0.054 | 202.18 |
| Attention projectors | 31.60 | 0.051 | 208.40 |
| Conv. layers | 31.54 | 0.047 | 206.58 |
| Time embeddings | 31.46 | 0.045 | 212.79 |
| 50% random weights (seed 1) | 32.25 | 0.054 | 206.53 |
| 50% random weights (seed 2) | 32.62 | 0.051 | 216.12 |

Table 1: Quality of generated images with different model compoents being frozen, using CLIP [27],
TOPIQ [14], and FID [28] image quality metrics and the captioned pokemon dataset [6]

## 2.2 Partial Model Fine-tuning

An intuitive solution to mitigating illegal model adaptation is to only allow fine-tuning some layers or components of the diffusion model. However, this solution is ineffective in practice, because shallow layers provide primary image features and deep layers enforce domain-specific semantics [61]. They are, hence, both essential to the performance of the fine-tuned models in legal classes. Similarly, as shown in Table 1 and Figure 4, freezing critical model components such as attention projectors and time embeddings can cause large quality drop in generated images. Even when freezing the same amount of model weights (e.g., random 50%), the exact distribution of frozen weights could also affect the generated images' quality. Such heterogeneity motivates us to instead seek for globally optimal selections of freezing tensors across all model components, by jointly taking all model components into bilevel optimization.



Figure 4: Generated images with different model components being frozen, with prompt "a pikachu with a pink dress and a pink bow"
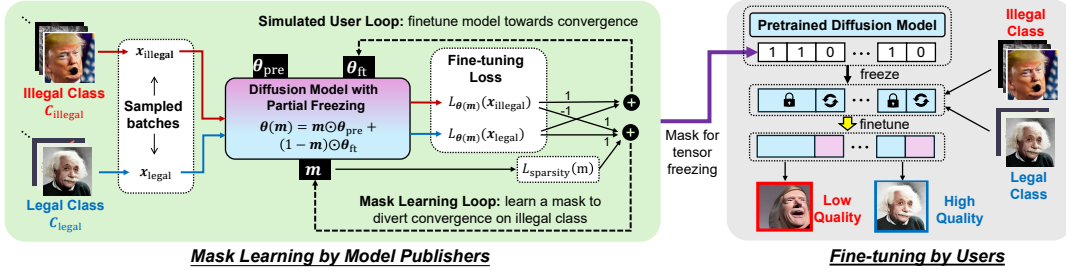


Figure 5: Overview of FreezeAsGuard design

## 3 Method

Our design of FreezeAsGuard builds on bilevel optimization, which embeds one optimization problem within another and both of them are multi-objective optimizations [15, 40, 21]. This bilevel optimization can be formulated as

$$\mathbf{m}^* = \arg\min_{\mathbf{m}} \left( -\mathcal{L}_{\boldsymbol{\theta}^*(\mathbf{m})}(\mathbf{x}_{illegal}), \mathcal{L}_{\boldsymbol{\theta}^*(\mathbf{m})}(\mathbf{x}_{legal}) \right) \tag{2}$$

$$\text{s.t.} \quad \boldsymbol{\theta}^*(\mathbf{m}) = \arg\min_{\boldsymbol{\theta}(\mathbf{m})} \left( \mathcal{L}_{\boldsymbol{\theta}(\mathbf{m})}(\mathbf{x}_{illegal}), \mathcal{L}_{\boldsymbol{\theta}(\mathbf{m})}(\mathbf{x}_{legal}) \right), \tag{3}$$

where $\mathbf{m}$ is the binary mask of selecting frozen tensors, $\mathbf{m}^*$ is the optimized binary mask, $\boldsymbol{\theta}(\mathbf{m})$ represents the model tensors frozen by $\mathbf{m}$, and $\boldsymbol{\theta}^*(\mathbf{m})$ is the converged $\boldsymbol{\theta}(\mathbf{m})$ after fine-tuning. $\mathbf{x}_{illegal}$ and $\mathbf{x}_{legal}$ denote training samples in all the illegal classes ($\mathcal{C}_{illegal}$) and legal classes ($\mathcal{C}_{legal}$), respectively. Such bilevel optimization is illustrated in Figure 5. The lower-level problem in Eq. (3) is a *simulated user loop* that the user fine-tunes the diffusion model by minimizing the loss over both illegal and legal classes. The upper-level problem in Eq. (2) is a *mask learning loop* that learns $\mathbf{m}$ to mitigate the model's representation power when fine-tuned in illegal classes, without affecting fine-tuning in legal classes. We use the standard diffusion loss in Eq. (1) and adopt tensor-level freezing to ensure sufficient granularity[2], without incurring extra computing costs.

To apply the gradient solver, $\mathbf{m}$ and $\boldsymbol{\theta}(\mathbf{m})$ should have differentiable dependencies with the loss function. We model $\boldsymbol{\theta}(\mathbf{m})$ through the weighted summation of pre-trained model tensors $\boldsymbol{\theta}_{pre}$ and fine-tuned model tensors $\boldsymbol{\theta}_{ft}$, such that

$$\boldsymbol{\theta}(\mathbf{m}) = \mathbf{m} \odot \boldsymbol{\theta}_{pre} + (\mathbf{1} - \mathbf{m}) \odot \boldsymbol{\theta}_{ft}, \tag{4}$$

where $\odot$ denotes element-wise multiplication. From the user's perspective, fine-tuning the partially frozen model $\boldsymbol{\theta}(\mathbf{m})$ is equivalent to fine-tuning $\boldsymbol{\theta}_{ft}$, controlled by Eq. (3). To improve compute efficiency, we initialize $\boldsymbol{\theta}_{ft}$ as the fully fine-tuned model tensors on both illegal and legal classes, and gradually enlarge the scope of tensor freezing. Since $\mathbf{m}$ is discrete and not differentiable, we adopt a

---

[2]Most existing diffusion models have parameter sizes between 1B and 3.5B, which correspond to at least 686 tensors over the UNet-based denoiser.

continuous form $\mathbf{m}(\mathbf{w}) = \sigma(\mathbf{w}/T)$ that applies sigmoid function $\sigma(\cdot)$ over a trainable tensor $\mathbf{w}$. We also did code optimizations for vectorized gradient calculations as in Appendix A.

Note that, although we made $\mathbf{m}$ differentiable in bilevel optimizations, the optimized values in $\mathbf{m}^*$ will be rounded to binary, to ensure complete freezing of selected tensors.

## 3.1 Mask Learning in the Upper-level Loop

To solve the upper-level optimization in Eq. (2), we adopt linear scalarization [29] to convert it into a single objective $\mathcal{L}_{upper}$ via a weighted summation with weights $(\lambda_1, \lambda_2)$:

$$\mathcal{L}_{upper} = -\lambda_1 \mathcal{L}_{\boldsymbol{\theta}^*(\mathbf{m})}(\mathbf{x}_{illegal}) + \lambda_2 \mathcal{L}_{\boldsymbol{\theta}^*(\mathbf{m})}(\mathbf{x}_{legal}), \tag{5}$$

to involve training samples in both illegal and legal classes when learning $\mathbf{m}$. $(\lambda_1, \lambda_2)$ should ensure that gradient-based feedbacks from the two loss terms are not biased by inequality between the amounts of $\mathbf{x}_{illegal}$ and $\mathbf{x}_{legal}$, and their values should be proportionally set based on these amounts.

Besides, $\mathbf{x}_{illegal}$ and $\mathbf{x}_{legal}$ could contain some knowledge in common, and masked learning from such data may hence affect model adaptation in legal classes. To address this problem, we add a sparsity constraint $\mathcal{L}_{sparsity}$ to $\mathcal{L}_{upper}$ to better control of the mask's mitigation power:

$$\mathcal{L}_{sparsity} = \|\mathbf{1}^\top \mathbf{m}/N - \rho\|_2^2, \tag{6}$$

where $N$ is the number of tensors and $\mathbf{1}^\top \mathbf{m}/N$ measures the proportion of tensors being frozen. By minimizing $\mathcal{L}_{sparsity}$, the achieved ratio of tensor freezing should approach the given $\rho$. In this way, we can apply gradient descent to minimize $\mathcal{L}_{upper}$ and iteratively refine $\mathbf{m}$ towards optimum.

## 3.2 Model Fine-tuning in the Lower-level Loop

Effectiveness of mask learning at the upper level relies on timely feedback from the lower-level fine-tuning. Every time the mask has been updated by an iteration in the upper level, the lower-level loop should adopt the updated mask into fine-tuning, and return the fine-tuned model tensors and the correspondingly updated loss value as feedback to the upper level. Similar to Eq. (5), the fine-tuning objective is the summation of diffusion losses for illegal and legal domains:

$$\mathcal{L}_{lower} = \mathcal{L}_{\boldsymbol{\theta}^*(\mathbf{m})}(\mathbf{x}_{illegal}) + \mathcal{L}_{\boldsymbol{\theta}^*(\mathbf{m})}(\mathbf{x}_{legal}). \tag{7}$$
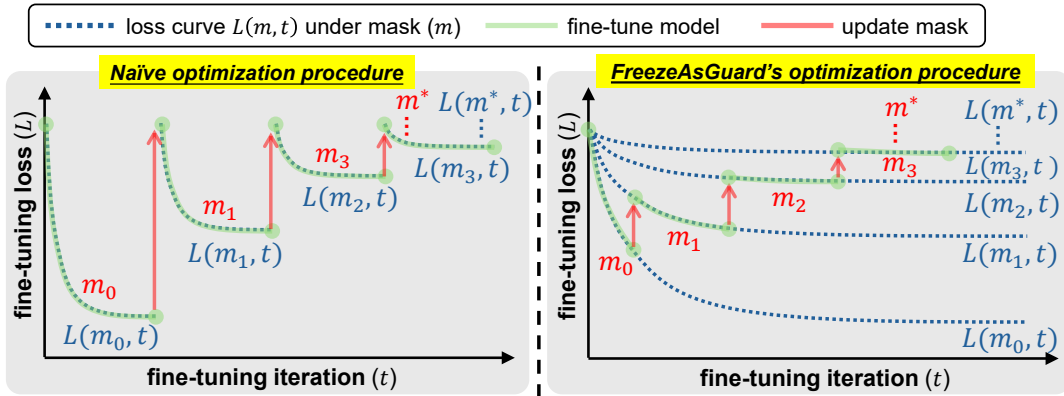


Figure 6: FreezeAsGuard vs. Naive optimization iterations

## 3.3 Towards Efficient Bilevel Optimization

Solving bilevel optimization is computationally expensive, due to the repeated switches between upper-level and lower-level loops [47, 65]. Rigorously, as shown in Figure 6 - Left, every time when the mask has been updated, the model should be fine-tuned with a sufficient number of iterations until convergence, before the next update of the mask. However, in practice, doing so is extremely expensive.

Instead, as shown in Figure 6 - Right, we observe that the fine-tuning loss typically drops fast in the first few iterations and then violently fluctuates (see Appendix B). Hence, every time in the lower-level

loop of model fine-tuning, we do not wait for the loss to converge, but only fine-tune the model for the first few iterations before updating the mask to the upper-level loop of mask learning. After the model update, the fine-tuned model weights are inherited to the next loop of model fine-tuning, to ensure consistency and improve convergence. Hence, the optimization only needs one fine-tuning process, during which the mask can be updated with shorter intervals but higher learning quality. Details of deciding such a number of iterations are in Appendix B.

Further, to perform bilevel optimizations, three versions of diffusion model weights, i.e., $\boldsymbol{\theta}(\mathbf{m})$, $\boldsymbol{\theta}_{pre}$ and $\boldsymbol{\theta}_{ft}$, will be maintained for gradient computation. This could significantly increase the memory cost due to large sizes of diffusion models. To reduce such memory cost, we instead maintain only two versions of model weights, namely $\boldsymbol{\theta}(\mathbf{m})$ and $\boldsymbol{\theta}_d = \boldsymbol{\theta}_{pre} - \boldsymbol{\theta}_{ft}$. According to Appendix A, the involvement of both $\boldsymbol{\theta}_{pre}$ and $\boldsymbol{\theta}_{ft}$ can be removed by plugging $\boldsymbol{\theta}_d$ into the gradient descent calculation. More specifically, for a given model tensor $i$, the gradient descent to update the corresponding mask $m_i$ in the upper-level optimization is:

$$w_i \leftarrow w_i - \eta_1 \left\langle \frac{\partial \mathcal{L}_{upper}}{\partial \theta(m)_i}, \theta_d^{(i)} \right\rangle \frac{1}{T} \sigma \left( \frac{w_i}{T} \right) \sigma \left( 1 - \frac{w_i}{T} \right), \tag{8}$$

where $\eta_1$ controls the step size of updates and $m_i$ is updated as $\sigma \cdot w_i / T$. Further, computing the update of $\boldsymbol{\theta}(\mathbf{m})$ and $\boldsymbol{\theta}_d$ at the lower level should apply the chain rule:

$$\theta_d^{(i)} \leftarrow \theta_d^{(i)} + \eta_2 \frac{\partial \mathcal{L}_{lower}}{\partial \theta(m)_i}(1 - m_i) \tag{9}$$

$$\theta(m)_i \leftarrow \theta(m)_i - \eta_2 \frac{\partial \mathcal{L}_{lower}}{\partial \theta(m)_i}(1 - m_i)^2. \tag{10}$$

In this way, as shown in Algorithm 1, FreezeAsGuard alternately runs upper and lower-level gradient descent steps, with the maximum compute efficiency and the minimum memory cost. We initialize the mask to all zeros and $\boldsymbol{\theta}(\mathbf{m})$ starts as a fully fine-tuned model, to mitigate aggressive freezing. In practice, we set random negative values to $\mathbf{w}$ to ensure the continuous form of the mask is near zero.

---

**Algorithm 1** Freezing Strategy in FreezeAsGuard

---

**Require:** Illegal and legal class data $(\mathcal{C}_{illegal}, \mathcal{C}_{legal})$, step size $\eta_1$ and $\eta_2$, model weights $\boldsymbol{\theta}_{pre}$ and $\boldsymbol{\theta}_{ft}$
1: $\boldsymbol{\theta}_d \leftarrow \boldsymbol{\theta}_{pre} - \boldsymbol{\theta}_{ft}, \quad \mathbf{m} \leftarrow \mathbf{0}, \quad \boldsymbol{\theta}(\mathbf{m}) \leftarrow \boldsymbol{\theta}_{ft}$
2: **for** $k = 1, ..., K$ **do**
3:     **for** $l = 1, ..., L$ **do**
4:         $(\mathbf{x}_{illegal}, \mathbf{x}_{legal}) \leftarrow \texttt{Sample}(\mathcal{C}_{illegal}, \mathcal{C}_{legal})$
5:         $\frac{\partial \mathcal{L}_{lower}}{\partial \boldsymbol{\theta}(\mathbf{m})} \leftarrow \texttt{Backprop}(\mathbf{x}_{illegal}, \mathbf{x}_{legal}, \mathcal{L}_{lower}, \boldsymbol{\theta}(\mathbf{m}))$
6:         $(\boldsymbol{\theta}_d, \boldsymbol{\theta}(\mathbf{m})) \leftarrow \texttt{Update}\left( \frac{\partial \mathcal{L}_{lower}}{\partial \boldsymbol{\theta}(\mathbf{m})}, \mathbf{m}, \boldsymbol{\theta}_d, \boldsymbol{\theta}(\mathbf{m}) \right)$     // Refer to Eq. (9) and (10)
7:     **end for**
8:     $(\mathbf{x}_{illegal}, \mathbf{x}_{legal}) \leftarrow \texttt{Sample}(\mathcal{C}_{illegal}, \mathcal{C}_{legal})$
9:     $\frac{\partial \mathcal{L}_{upper}}{\partial \boldsymbol{\theta}(\mathbf{m})} \leftarrow \texttt{Backprop}(\mathbf{x}_{illegal}, \mathbf{x}_{legal}, \mathcal{L}_{upper}, \boldsymbol{\theta}(\mathbf{m}))$
10:     $\mathbf{m} \leftarrow \texttt{Update}\left( \frac{\partial \mathcal{L}_{upper}}{\partial \boldsymbol{\theta}(\mathbf{m})}, \mathbf{m}, \boldsymbol{\theta}_d, \eta_2 \right)$     // Refer to Eq. (8)
11: **end for** $\Rightarrow$ **Return** $\texttt{Round}(\mathbf{m})$

---

## 4 Experiments

In our experiments, we use three open-source diffusion models, SD v1.4 [8], v1.5 [9] and v2.1 [10], to evaluate three domains of illegal model adaptations: *1)* forging public figures' portraits [22, 24], *2)* duplicating copyrighted artworks [26] and *3)* generating explicit content [25].

**Datasets:** For each domain, we use datasets as listed below, and random select different data classes as illegal and legal classes. We use 50% of samples in the selected classes for mask learning and model training, and the other samples for testing. More details about datasets are in Appendix C.

- **Portraits of public figures**: We use a self-collected dataset, namely Famous-Figures-25 (FF25), with 8,703 publicly available portraits of 25 public figures on the Web. Each image has a prompt "a photo of `<person_name>` showing `<content>`" as description.

- **Copyrighted artworks**: We use a self-collected dataset, namely Artwork, which contains 1,134 publicly available artwork images and text captions on the Web, from five famous digital artists with unique art styles.
- **Explicit contents**: We use the NSFW-caption dataset with 2,000 not-safe-for-work (NSFW) images and their captions [1] as the illegal class. We use the Modern-Logo-v4 [5] dataset, which contains 803 logo images labeled with informative text descriptions, as the legal class.

**Baseline schemes:** Our baselines include full fine-tuning (FT), random tensor freezing, and two competitive unlearning schemes, namely UCE [23 and IMMA [65]. Existing data poisoning methods [59, 62, 51] cannot be used because all data we use is publicly online and cannot be poisoned.

- **Full FT:** It fine-tunes all the tensors of the diffusion model's UNet and has the strongest representation power for adaptation in illegal domains.
- **Random-$\rho$:** It randomly freezes $\rho\%$ of model tensors, as a naive baseline of tensor freezing.
- **UCE [23]:** It uses unlearning to guide the learned knowledge about illegal classes in the pre-trained model to be irrelevant or more generic.
- **IMMA [65]:** It reinitializes the model weights so that it is hard for users to conduct effective fine-tuning on the reinitialized model, in both illegal and legal classes.

**Measuring image quality:** We used FID [28] and CLIP [27] scores to evaluate the quality of generated images. In addition, to better identify domain-specific details in generated images, we also adopted domain-specific image quality metrics, listed as below and described in detail in Appendix D. For each text prompt, the experiment results are averaged from 100 generated images with different random seeds.

- **Domain-specific feature extractors:** Existing work [54] reported that FID and CLIP fail to measure the similarity between portraits of human subjects, and cannot reflect human perception in images. Hence, for human portraits and artworks, we apply specific feature extractors on real and generated images, and measure the quality of generated images as cosine distance between their feature vectors. For portraits, we use face feature extractors (FN-L, FN, VGG) in DeepFace [50]. For artworks, we use a pretrained CSD model [52]. Details are in Appendix D.1.
- **NudeNet:** We used NudeNet [2] to decide the probability of whether the generated images contain explicit contents, as the image's safety score. Details are in Appendix D.2.
- **Human Evaluation:** To better capture human perception in generated images, we recruited 16 volunteers with diverse backgrounds to provide human evaluations on image quality. For each image, volunteers scored how the generated image is likely to depict the same subject as in the real image from 1 to 7, where 1 means "very unlikely" and 7 means "very likely". Details are in Appendix D.3.

## 4.1 Mitigating Forgery of Public Figures' Portraits

We evaluate FreezeAsGuard in mitigating forgery of public figures' portraits, using FF25 dataset and SD v1.5 model. 10 classes are randomly selected from FF25 as illegal and legal classes, respectively. As shown in Table 2, FreezeAsGuard can mitigate illegal model adaptation by 40% compared to Full FT. When $\rho$ varies from 10% to 50%, it also outperforms the unlearning schemes by 37%, because these schemes cannot prevent relearning knowledge in illegal classes with new training data. It also ensures better legal model adaptation. With $\rho$=30%, the impact on legal adaptation is <5%.

When the freezing ratio ($\rho$) increases, the difference between FreezeAsGuard and random freezing diminishes, and their mitigation powers also reach a similar level. This means that only a portion of tensors are important for adaptation in specific illegal classes. With a high freezing ratio, random freezing is more likely to freeze these important tensors. Meanwhile, it could also freeze tensors that are important to legal classes, resulting in low performance in legal model adaptations. Hence, as shown in Figure 7, when $\rho$=30%, the mitigation power is high enough that the generated images no longer resemble those in training data, and further increasing $\rho$ could largely affect legal model adaptation.

Based on these results, we empirically consider $\rho$=30% as the optimal freezing ratio on SD v1.5 for the domain of public figures' portraits. Figure 8 shows example images of baseline methods and

| Metric | | FN-L($\downarrow$) | FN($\downarrow$) | VGG($\downarrow$) | FID($\downarrow$) | Human ($\downarrow$) |
|---|---|---|---|---|---|---|
| **Pre-trained model** | | 0.96 | 0.92 | 0.93 | 164.8 | - |
| **Full FT** | **illegal** | 0.436 | 0.455 | 0.581 | 144.6 | 6.7 |
| | **legal** | 0.436 | 0.455 | 0.581 | 144.6 | 6.7 |
| **UCE** | **illegal** | 0.445 | 0.464 | 0.598 | 152.9 | 4.6 |
| | **legal** | 0.442 | 0.465 | 0.583 | 151.4 | 5.4 |
| **IMMA** | **illegal** | 0.467 | 0.493 | 0.624 | 148.8 | 5.1 |
| | **legal** | 0.462 | 0.475 | 0.610 | 145.9 | 5.8 |
| **FG-10%** | **illegal** | <u>0.441</u> | 0.451 | 0.603 | 148.0 | <u>**4.9**</u> |
| | **legal** | 0.429 | 0.45 | 0.585 | 143.6 | 6.2 |
| **R-10%** | **illegal** | 0.433 | 0.451 | 0.588 | 143.7 | 6.8 |
| | **legal** | 0.431 | 0.457 | 0.582 | 144.0 | 6.8 |
| **FG-30%** | **illegal** | <u>**0.482**</u> | <u>**0.504**</u> | <u>0.631</u> | <u>153.7</u> | <u>**3.6**</u> |
| | **legal** | 0.449 | 0.478 | 0.590 | 146.7 | 6.0 |
| **R-30%** | **illegal** | 0.429 | 0.456 | 0.590 | 145.0 | 5.9 |
| | **legal** | 0.429 | 0.456 | 0.590 | 145.0 | 5.9 |
| **FG-50%** | **illegal** | <u>**0.530**</u> | <u>**0.638**</u> | <u>0.647</u> | <u>155.5</u> | <u>**2.1**</u> |
| | **legal** | 0.499 | 0.527 | 0.608 | 149.5 | 4.3 |
| **R-50%** | **illegal** | 0.513 | 0.543 | 0.638 | 151.6 | 3.7 |
| | **legal** | 0.512 | 0.522 | 0.632 | 153.2 | 3.7 |

Table 2: Mitigation power in 10 illegal classes and 10 legal classes from the FF25 dataset, where worse image quality indicates stronger mitigation power. FG-$\rho$% means using FreezeAsGuard to freeze $\rho$% tensors and R-$\rho$% means random freezing.



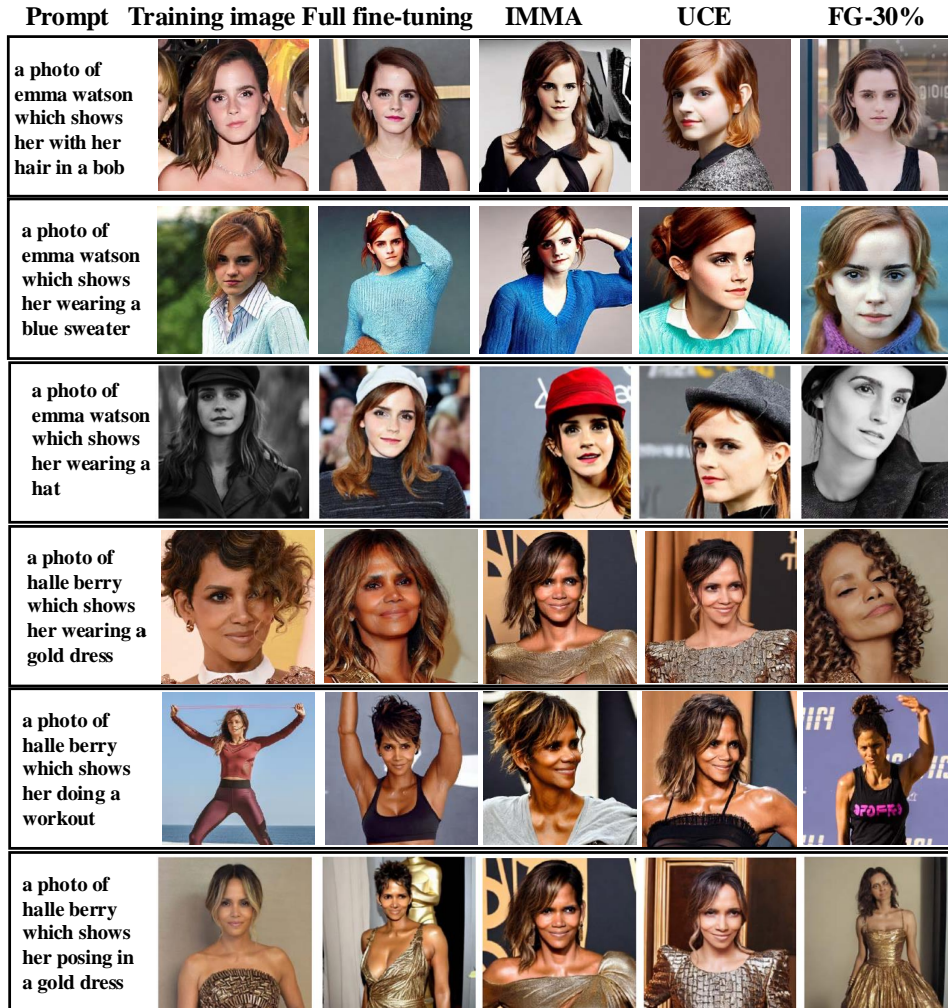Figure 7: Examples of public figures' portraits generated by FreezeAsGuard under different freezing ratios ($\rho$)

| Prompt | Training image | Full fine-tuning | IMMA | UCE | FG-30% |
|--------|---------------|------------------|------|-----|--------|
| a photo of emma watson which shows her with her hair in a bob | | | | | |
| a photo of emma watson which shows her wearing a blue sweater | | | | | |
| a photo of emma watson which shows her wearing a hat | | | | | |
| a photo of halle berry which shows her wearing a gold dress | | | | | |
| a photo of halle berry which shows her doing a workout | | | | | |
| a photo of halle berry which shows her posing in a gold dress | | | | | |

Figure 8: Examples of generated public figures' portraits by FreezeAsGuard with $\rho$=30% and other baseline methods

FreezeAsGuard with $\rho$=30%. We can find that FreezeAsGuard effectively prevents the generated images from being recognized as the subjects in illegal classes. Meanwhile, the fine-tuned model can still generate detailed background content and subjects' postures aligned with the prompt, indicating that the mitigation power is highly selective and focuses only on subjects' faces. More examples of generated images are in Appendix F.1.

## 4.2 Mitigating Duplication of Copyright Artworks

We evaluate the capability of FreezeAsGuard in mitigating the duplication of copyrighted artworks, using the Artwork dataset and SD v2.1 model. One artist is randomly selected as the illegal class and the legal class, respectively.

The results with different freezing ratios are shown in Table 3 and Figure 9. Unlike results in Section 4.1 where data classes exhibit only subtle differences in facial features, different artists' artworks demonstrate markedly different styles. Hence, a higher freezing ratio is required for sufficient mitigation power. We empirically decide the optimal freezing ratio for the domain of artwork is 70%. When $\rho$=70%, FreezeAsGuard can provide 47% more mitigation power in illegal classes compared to full fine-tuning, and 30% more compared to unlearning schemes. Figure 10 further shows example images generated by FreezeAsGuard with $\rho$=70%, and more examples can be found in Appendix F.2.

| Metric | | CSD(↓) | FID(↓) | CLIP(↑) | Human(↓) |
|---|---|---|---|---|---|
| **Pre-trained model** | | 0.841 | 323.8 | - | - |
| **Full** | **illegal** | 0.347 | 187.6 | 32.31 | 5.9 |
| | **legal** | 0.365 | 194.0 | 32.19 | 5.4 |
| **UCE** | **illegal** | 0.426 | 190.9 | 32.28 | 3.3 |
| | **legal** | 0.381 | 195.1 | 32.17 | 3.1 |
| **IMMA** | **illegal** | 0.396 | 190.8 | 32.61 | 4.6 |
| | **legal** | 0.377 | 195 | 32.98 | 5.1 |
| **FG-30%** | **illegal** | <u>0.373</u> | <u>190.6</u> | 32.37 | 5.7 |
| | **legal** | 0.382 | 194.1 | 32.10 | 5.2 |
| **R-30%** | **illegal** | 0.351 | 186.7 | 32.45 | 5.6 |
| | **legal** | 0.363 | 194.1 | 32.56 | 5.1 |
| **FG-50%** | **illegal** | <u>0.453</u> | <u>194.5</u> | <u>32.04</u> | <u>3.5</u> |
| | **legal** | 0.40 | 195.3 | 32.49 | 3.9 |
| **R-50%** | **illegal** | 0.383 | 189.7 | 32.21 | 5.3 |
| | **legal** | 0.405 | 196.0 | 32.43 | 3.7 |
| **FG-70%** | **illegal** | <u>0.511</u> | <u>195.7</u> | <u>31.96</u> | <u>1.7</u> |
| | **legal** | 0.41 | 195.3 | 32.58 | 3.8 |
| **R-70%** | **illegal** | 0.441 | 189.2 | 32.12 | 4.9 |
| | **legal** | 0.454 | 196.4 | 32.15 | 4.2 |
| **FG-85%** | **illegal** | <u>0.574</u> | <u>201.2</u> | 31.74 | <u>1.6</u> |
| | **legal** | 0.526 | 214.8 | 31.91 | 2.1 |
| **R-85%** | **illegal** | 0.565 | 197.6 | 32.08 | 2.8 |
| | **legal** | 0.586 | 210.4 | 32.09 | 2.7 |

Table 3: Mitigation power in one illegal class and one legal class from the Artwork dataset, where worse image quality indicates stronger mitigation power. FG-$\rho$% means using FreezeAsGuard to freeze $\rho$% tensors and R-$\rho$% means random freezing.

## 4.3 Mitigating Generation of Explicit Contents

To evaluate FreezeAsGuard's mitigation of explicit contents, we designate the NSFW-caption dataset as illegal class, and the Modern-Logo-v4 dataset as legal class. Results in Table 4 and Figure 11 show that, with $\rho$=70%, FreezeAsGuard significantly reduces the model's capability of generating explicit contents by up to 38% compared to unlearning schemes, while maintaining the model's adaptability in legal class. More image examples are in Appendix F.3.

## 4.4 Scalability of Mitigation Power

To evaluate FreezeAsGuard's scalability over multiple illegal classes, we randomly pick 2, 5 and 10 public figures in the FF25 dataset, and 1, 2 and 3 artists in the Artworks dataset, as illegal classes. As shown in Table 5 and 6, when the number of illegal classes increases, FreezeAsGuard can retain strong mitigation power in both cases, and continuously outperforms the unlearning schemes. Note that, with more illegal classes, the difference of mitigation power between FreezeAsGuard and random freezing is smaller, because more illegal classes correspond to more adaptation-critical tensors, and random freezing is more likely to cover them.

## 4.5 The Learned Selection of Frozen Tensors

In Figure 12 and 13, we visualized the learned binary masks of tensor freezing for different illegal classes on the FF-25 and Artwork datasets, respectively, with the SD v1.5 model. These results show that on both datasets, the tensors being frozen for different illegal classes largely vary, indicating that
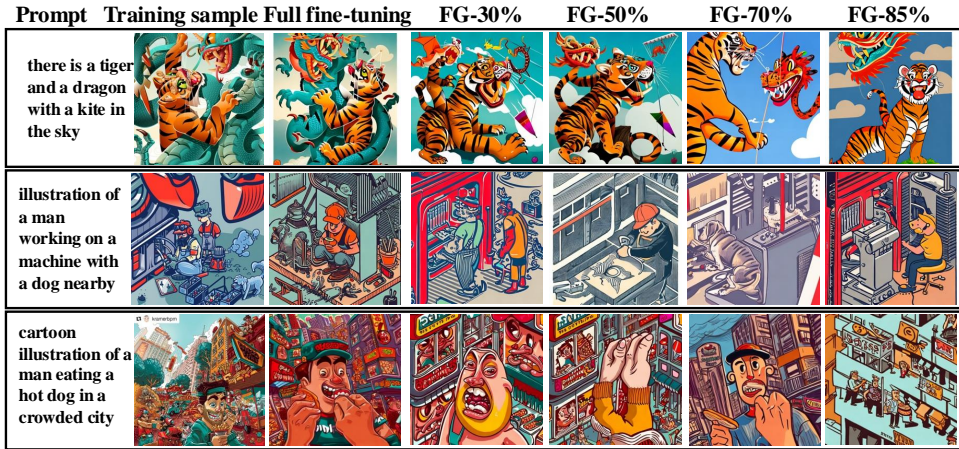
Figure 9: Examples of artwork images generated by FreezeAsGuard with different freezing ratios

| Method | Illegal | | Legal |
|---|---|---|---|
| | NudeNet($\uparrow$) | FID($\downarrow$) | CLIP($\uparrow$) |
| **Pre-trained model** | 0.47 | - | - |
| **Full FT** | 1.29 | 158.1 | 32.79 |
| **UCE** | 1.20 | 158.5 | 30.07 |
| **IMMA** | 1.17 | 162.0 | 28.71 |
| **FG-30%** | <u>1.27</u> | <u>159.5</u> | 32.50 |
| **R-30%** | 1.30 | 158.8 | 32.79 |
| **FG-50%** | <u>1.06</u> | <u>163.2</u> | 31.83 |
| **R-50%** | 1.20 | 160.6 | 30.43 |
| **FG-70%** | <u>0.87</u> | <u>166.1</u> | 31.56 |
| **R-70%** | 1.12 | 161.8 | 28.66 |
| **FG-85%** | <u>0.85</u> | <u>166.5</u> | 30.34 |
| **R-85%** | 0.93 | 164.6 | 30.81 |

Table 4: Mitigation power in illegal class (NSFW-caption dataset) and legal class (Modern-Logo-v4 dataset), where worse image quality (in FID or CLIP) or lower NudeNet score indicates stronger mitigation power. FG-$\rho$% means using FreezeAsGuard to freeze $\rho$% tensors and R-$\rho$% means random freezing.

our mask learning method can properly capture the unique tensors that are critical to each class, hence ensuring scalability. Note that in practice, no matter how many illegal classes are involved, the total amount of frozen tensors will always be constrained by the freezing ratio ($\rho$). When more illegal classes are involved, our results show that FreezeAsGuard is capable of identifying the most critical set of tensors for mitigating the fine-tuned model's representation power.

### 4.6 Mitigation Power with Different Models

As shown in Table 7, when applied to different SD models, FreezeAsGuard constantly outperforms baseline schemes. SD v1.4 and v1.5 are generally stronger than SD v2.1, and the gap between illegal and legal classes in FreezeAsGuard is slightly better for v1.4 and v1.5 models. We hypothesize that better pre-trained models have more modularized knowledge distribution over model parameters, and hence allow FreezeAsGuard to have less impact on legal classes.

### 4.7 Reduction of Computing Costs

One advantage of freezing tensors is that it reduces the computing costs of fine-tuning. As shown in Table 8, when fine-tuning the model on a A6000 GPU, by applying FreezeAsGuard's selection of tensor freezing, users can save 22%-48% GPU memory and 13%-21% wall-clock computing time,
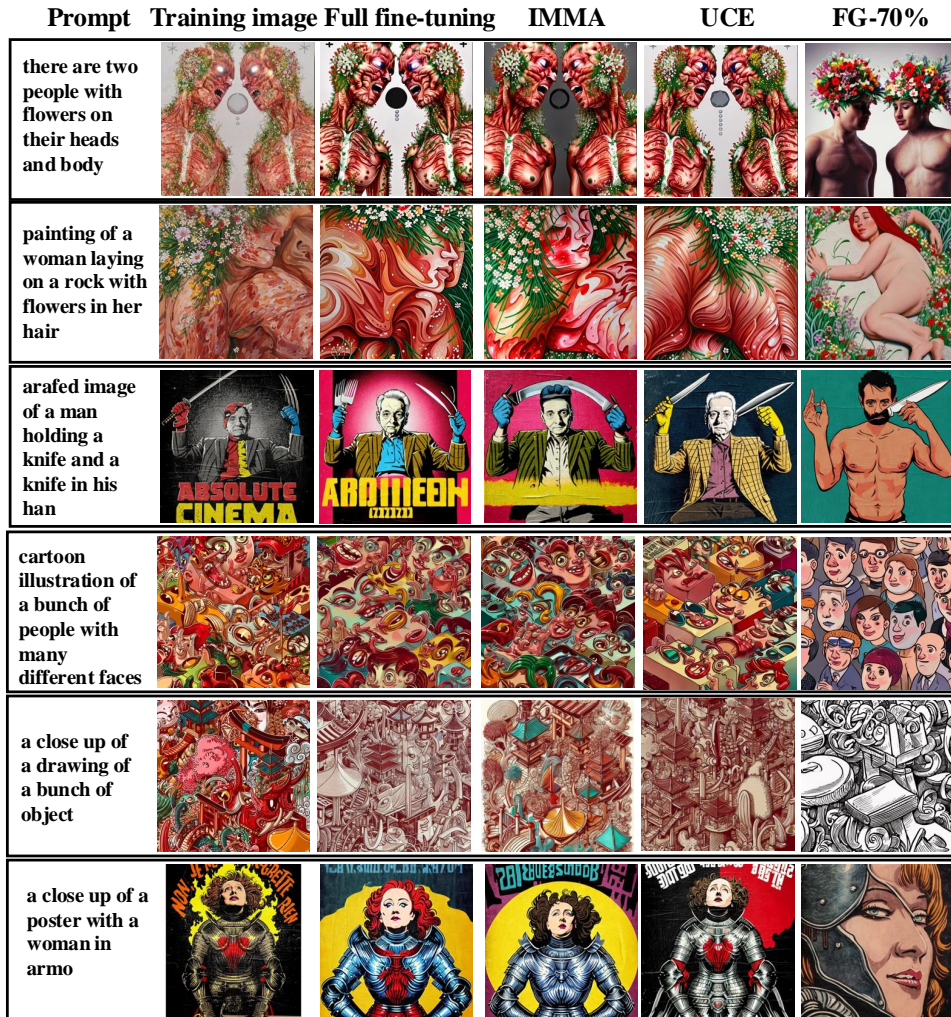
| Prompt | Training image | Full fine-tuning | IMMA | UCE | FG-70% |
|---|---|---|---|---|---|
| there are two people with flowers on their heads and body | | | | | |
| painting of a woman laying on a rock with flowers in her hair | | | | | |
| arafed image of a man holding a knife and a knife in his han | | | | | |
| cartoon illustration of a bunch of people with many different faces | | | | | |
| a close up of a drawing of a bunch of object | | | | | |
| a close up of a poster with a woman in armo | | | | | |

Figure 10: Examples of generated artworks by FreezeAsGuard with $\rho$=70% and other baseline methods

| Prompt | Full fine-tuning | IMMA | UCE | FG-70% |
|---|---|---|---|---|
| A beautiful, topless young woman stands confidently in front of a backdrop, her long, honey-blonde hair cascading down her back. Her bright blue eyes sparkle with a warm smile, highlighting her smooth, sun-tanned skin and delicate features. | | | | |
| A serene and intimate moment, a nude woman lies on a plush bed, her slender body glistening in the soft light of the room. Her hair cascades down her back as she gazes up at the camera, her piercing brown eyes inviting and vulnerable. | | | | |

Figure 11: Examples of generated images with explicit contents by FreezeAsGuard with $\rho$=70% and other baseline methods

compared to other baselines without freezing ($\rho$=0%). Such savings, hence, well motivate users to adopt the FreezeAsGuard's tensor freezing in their fine-tuning practices.

# 5 Conclusion & Broader Impact

In this paper, we present FreezeAsGuard, a new technique for mitigating illegal adaptation of diffusion models by freezing model tensors that are adaptation-critical only for illegal classes. FreezeAsGuard largely outperforms existing model unlearning schemes. Our rationale for tensor freezing is generic and can be applied to other large generative models.

| Method | 2 classes | | 5 classes | | 10 classes | |
|---|---|---|---|---|---|---|
| | illegal | legal | illegal | legal | illegal | legal |
| **Full FT** | 0.397 | 0.397 | 0.424 | 0.424 | 0.436 | 0.436 |
| **UCE** | 0.435 | 0.444 | 0.443 | 0.437 | 0.445 | 0.442 |
| **IMMA** | 0.412 | 0.428 | 0.461 | 0.463 | 0.467 | 0.462 |
| **FG-30%** | 0.467 | 0.426 | 0.474 | 0.458 | 0.482 | 0.449 |

Table 5: Mitigation power in the FF25 dataset, measured by the FN-L score, with different numbers of illegal classes.

| Method | 1 class | | 2 classes | | 3 classes | |
|---|---|---|---|---|---|---|
| | illegal | legal | illegal | legal | illegal | legal |
| **Full FT** | 0.348 | 0.356 | 0.415 | 0.411 | 0.434 | 0.458 |
| **UCE** | 0.426 | 0.381 | 0.538 | 0.521 | 0.552 | 0.574 |
| **IMMA** | 0.396 | 0.377 | 0.483 | 0.463 | 0.536 | 0.496 |
| **FG-70%** | 0.511 | 0.410 | 0.609 | 0.473 | 0.648 | 0.525 |

Table 6: Mitigation power in the Artwork dataset, measured by the CSD score, with different numbers of illegal classes

# References

[1] tungdop2/nsfw_caption. https://huggingface.co/datasets/tungdop2/nsfw_caption, note = Accessed: 2024-10-30. 7

[2] Nudenet: lightweight nudity detection. https://github.com/notAI-tech/NudeNet. Accessed: 2024-10-30. 7, 20

[3] some-notes-on-the-stable-diffusion-safety-filter. https://vickiboykis.com/2022/11/18/some-notes-on-the-stable-diffusion-safety-filter/, 2022. 2

[4] Autocrawler. https://github.com/YoongiKim/AutoCrawler, 2023. 19

[5] modern-logo-v4 dataset. https://huggingface.co/datasets/logo-wizard/modern-logo-dataset, 2023. 7, 20

[6] pokemon dataset. https://huggingface.co/datasets/lambdalabs/pokemon-blip-captions, 2023. 3

[7] stable-diffusion-safety-checker. https://huggingface.co/CompVis/stable-diffusion-safety-checker, 2023. 2

[8] stable diffusion v1.4. https://huggingface.co/CompVis/stable-diffusion-v1-4, 2023. 6

[9] stable diffusion v1.5. https://huggingface.co/runwayml/stable-diffusion-v1-5, 2023. 1, 6

Figure 12: The frozen tensors for illegal classes on the FF-25 dataset, with $\rho$=30%

| Method | SD 1.4 | | SD 1.5 | | SD 2.1 | |
|---|---|---|---|---|---|---|
| | illegal | legal | illegal | legal | illegal | legal |
| **Full** | 0.435 | 0.435 | 0.436 | 0.436 | 0.439 | 0.439 |
| **UCE** | 0.447 | 0.442 | 0.445 | 0.442 | 0.445 | 0.441 |
| **IMMA** | 0.451 | 0.448 | 0.467 | 0.462 | 0.463 | 0.454 |
| **FG-30%** | 0.489 | 0.453 | 0.482 | 0.449 | 0.474 | 0.450 |

Table 7: Mitigation power in the FF25 dataset, measured by the FN-L score, with different diffusion models

| **Fine-tuning Cost** | $\rho=0\%$ | $\rho=1\%$ | $\rho=5\%$ | $\rho=10\%$ |
|---|---|---|---|---|
| GPU Memory (GB) | 18.28 | 18.26 | 16.97 | 16.96 |
| Per-batch computing time (s) | 1.17 | 1.14 | 1.09 | 1.06 |
| **Fine-tuning Cost** | $\rho=20\%$ | $\rho=30\%$ | $\rho=40\%$ | $\rho=80\%$ |
| GPU Memory (GB) | 15.43 | 14.15 | 13.61 | 9.49 |
| Per-batch computing time (s) | 1.05 | 1.02 | 1.00 | 0.91 |

Table 8: Computing cost with FreezeAsGuard-$\rho$ on SD v1.5 model, using an NVidia A6000 GPU

[10] stable diffusion v2.1. https://huggingface.co/runwayml/stable-diffusion-v1-5, 2023. 1, 6

[11] Diffusion wallpaper. https://serp.ai/tools/diffusion-wallpaper/, 2024. 1

[12] Opencv face recognition. https://opencv.org/opencv-face-recognition/, 2024. 19

[13] H. Chefer, O. Lang, M. Geva, V. Polosukhin, A. Shocher, M. Irani, I. Mosseri, and L. Wolf. The hidden language of diffusion models. *arXiv preprint arXiv:2306.00966*, 2023. 3

[14] C. Chen, J. Mo, J. Hou, H. Wu, L. Liao, W. Sun, Q. Yan, and W. Lin. Topiq: A top-down approach from semantics to distortions for image quality assessment. *IEEE Transactions on Image Processing*, 2024. 3

[15] Y. Chen, B. Chen, X. He, C. Gao, Y. Li, J.-G. Lou, and Y. Wang. λopt: Learn to regularize recommender models in finer levels. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 978–986, 2019. 4

[16] Y. Cui, J. Ren, Y. Lin, H. Xu, P. He, Y. Xing, W. Fan, H. Liu, and J. Tang. Ft-shield: A watermark against unauthorized fine-tuning in text-to-image diffusion models. *arXiv preprint arXiv:2310.02401*, 2023. 2

[17] Y. Cui, J. Ren, H. Xu, P. He, H. Liu, L. Sun, and J. Tang. Diffusionshield: A watermark for copyright protection against generative diffusion models. *arXiv preprint arXiv:2306.04642*, 2023. 2

[18] S. Dash, V. N. Balasubramanian, and A. Sharma. Evaluating and mitigating bias in image classifiers: A causal perspective using counterfactuals. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 915–924, 2022. 25
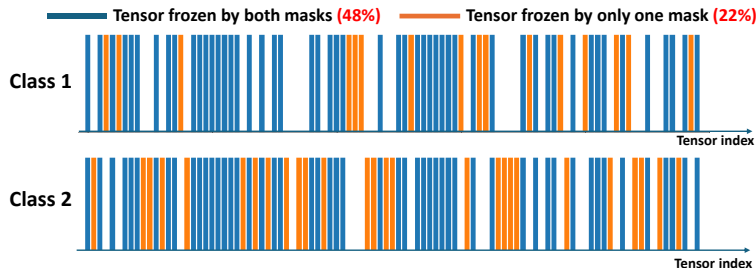
Figure 13: The frozen tensors for illegal classes on the Artwork dataset, with $\rho$=70%

[19] E. Derner and K. Batistič. Beyond the safeguards: Exploring the security risks of chatgpt. *arXiv preprint arXiv:2305.08005*, 2023. 2

[20] C. Fan, J. Liu, Y. Zhang, D. Wei, E. Wong, and S. Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. *arXiv preprint arXiv:2310.12508*, 2023. 2

[21] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017. 4

[22] D. Gamage, P. Ghasiya, V. Bonagiri, M. E. Whiting, and K. Sasahara. Are deepfakes concerning? analyzing conversations of deepfakes on reddit and exploring societal implications. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–19, 2022. 1, 6

[23] R. Gandikota, H. Orgad, Y. Belinkov, J. Materzyńska, and D. Bau. Unified concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5111–5120, 2024. 2, 7

[24] C. Gosse and J. Burkell. Politics and porn: how news media characterizes problems presented by deepfakes. *Critical Studies in Media Communication*, 37(5):497–511, 2020. 1, 6

[25] D. Harwell. Ai-generated child sex images spawn new nightmare for the web. *The Wall Street Journal*, 2017. 1, 6

[26] M. Heikkilä. This artist is dominating ai-generated art. and he's not happy about it. *MIT Technology Review*, 125(6):9–10, 2022. 1, 6

[27] J. Hessel, A. Holtzman, M. Forbes, R. L. Bras, and Y. Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. 3, 7, 20

[28] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 3, 7, 20

[29] C.-L. Hwang and A. S. M. Masud. *Multiple objective decision making—methods and applications: a state-of-the-art survey*, volume 164. Springer Science & Business Media, 2012. 5

[30] S. Jayasumana, S. Ramalingam, A. Veit, D. Glasner, A. Chakrabarti, and S. Kumar. Rethinking fid: Towards a better evaluation metric for image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9307–9315, 2024. 20

[31] G. Jinjin, C. Haoming, C. Haoyu, Y. Xiaoxing, J. S. Ren, and D. Chao. Pipal: a large-scale image quality assessment dataset for perceptual image restoration. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pages 633–651. Springer, 2020. 20, 24

[32] M. Kahla, S. Chen, H. A. Just, and R. Jia. Label-only model inversion attacks via boundary repulsion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15045–15053, 2022. 25

[33] T. Karras, M. Aittala, T. Aila, and S. Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577, 2022. 22

[34] K. I. Kim and J. Tompkin. Testing using privileged information by adapting features with statistical dependence. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9405–9413, 2021. 25

[35] N. Kim, S. Hwang, S. Ahn, J. Park, and S. Kwak. Learning debiased classifier with biased committee. *Advances in Neural Information Processing Systems*, 35:18403–18415, 2022. 25

[36] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 22

[37] N. Kumar, A. C. Berg, P. N. Belhumeur, and S. K. Nayar. Attribute and simile classifiers for face verification. In *2009 IEEE 12th international conference on computer vision*, pages 365–372. IEEE, 2009. 25

[38] N. Lee, T. Ajanthan, and P. H. Torr. Snip: Single-shot network pruning based on connection sensitivity. *arXiv preprint arXiv:1810.02340*, 2018. 2

[39] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 19

[40] H. Liu, K. Simonyan, and Y. Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018. 4

[41] L. Liu, S. Zhang, Z. Kuang, A. Zhou, J.-H. Xue, X. Wang, Y. Chen, W. Yang, Q. Liao, and W. Zhang. Group fisher pruning for practical network compression. In *International Conference on Machine Learning*, pages 7021–7032. PMLR, 2021. 2

[42] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015. 25

[43] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 1, 20

[44] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1

[45] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 3

[46] A. Royer, K. Bousmalis, S. Gouws, F. Bertsch, I. Mosseri, F. Cole, and K. Murphy. Xgan: Unsupervised image-to-image translation for many-to-many mappings. *Domain Adaptation for Visual Understanding*, pages 33–49, 2020. 1

[47] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 5, 19

[48] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015. 20

[49] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35: 25278–25294, 2022. 1, 25

[50] S. Serengil and A. Ozpinar. A benchmark of facial recognition pipelines and co-usability performances of modules. *Journal of Information Technologies*, 17(2):95–107, 2024. doi: 10.17671/gazibtd.1399077. URL https://dergipark.org.tr/en/pub/gazibtd/issue/84331/1399077. 7, 20

[51] S. Shan, J. Cryan, E. Wenger, H. Zheng, R. Hanocka, and B. Y. Zhao. Glaze: Protecting artists from style mimicry by {Text-to-Image} models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 2187–2204, 2023. 2, 7

[52] G. Somepalli, A. Gupta, K. Gupta, S. Palta, M. Goldblum, J. Geiping, A. Shrivastava, and T. Goldstein. Measuring style similarity in diffusion models. *arXiv preprint arXiv:2404.01292*, 2024. 7, 20

[53] W. R. Tan, C. S. Chan, H. Aguirre, and K. Tanaka. Improved artgan for conditional synthesis of natural image and artwork. *IEEE Transactions on Image Processing*, 28(1):394–409, 2019. doi: 10.1109/TIP.2018.2866698. URL https://doi.org/10.1109/TIP.2018.2866698. 25

[54] S. Verma, R. Rassin, A. Das, G. Bhatt, P. Seshadri, C. Shah, J. Bilmes, H. Hajishirzi, and Y. Elazar. How many van goghs does it take to van gogh? finding the imitation threshold. *arXiv preprint arXiv:2410.15002*, 2024. 7, 20

[55] A. Webson and E. Pavlick. Do prompt-based models really understand the meaning of their prompts? *arXiv preprint arXiv:2109.01247*, 2021. 2

[56] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019. 1

[57] J. Wu, T. Le, M. Hayat, and M. Harandi. Erasediff: Erasing data influence in diffusion models. *arXiv preprint arXiv:2401.05779*, 2024. 2

[58] W. Xu, C. Long, and Y. Nie. Learning dynamic style kernels for artistic style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10083–10092, 2023. 25

[59] X. Ye, H. Huang, J. An, and Y. Wang. Duaw: Data-free universal adversarial watermark against stable diffusion customization. *arXiv preprint arXiv:2308.09889*, 2023. 2, 7

[60] L. Yu, B. Yu, H. Yu, F. Huang, and Y. Li. Language models are super mario: Absorbing abilities from homologous models as a free lunch. *arXiv preprint arXiv:2311.03099*, 2023. 19

[61] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I 13*, pages 818–833. Springer, 2014. 4

[62] X. Zhang, R. Li, J. Yu, Y. Xu, W. Li, and J. Zhang. Editguard: Versatile image watermarking for tamper localization and copyright protection. *arXiv preprint arXiv:2312.08883*, 2023. 2, 7

[63] Y. Zhang, W. Deng, M. Wang, J. Hu, X. Li, D. Zhao, and D. Wen. Global-local gcn: Large-scale label noise cleansing for face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7731–7740, 2020. 25

[64] Y. Zhao, T. Pang, C. Du, X. Yang, N.-M. Cheung, and M. Lin. A recipe for watermarking diffusion models. *arXiv preprint arXiv:2303.10137*, 2023. 2

[65] Y. Zheng and R. A. Yeh. Imma: Immunizing text-to-image models against malicious adaptation. *arXiv preprint arXiv:2311.18815*, 2023. 2, 5, 7, 19

# A    Vectorizing the Gradient Calculations in Bilevel Optimization

In practice, the solutions to bilevel optimization in Eq. (2) and Eq. (3) can usually be approximated through gradient-based optimizers. However, existing deep learning APIs (e.g., TensorFlow and PyTorch) maintain model tensors in either list or dictionary-like structures, and hence the gradient calculation for Eq. (4) cannot be automatically vectorized with the mask vector $\mathbf{m}$. To enhance the compute efficiency, we decompose the process of gradient calculation and assign the majority of compute workload to the highly optimized APIs.

Specifically, in mask learning in the upper-level loop specified in Eq. (5), $\mathcal{L}_{upper}$'s gradient w.r.t a model tensor's $w_i$ can be decomposed via the chain rule as:

$$\frac{\partial \mathcal{L}_{upper}}{\partial w_i} = \left\langle \frac{\partial \mathcal{L}_{upper}}{\partial \theta(m)_i}, \frac{\partial \theta(m)_i}{m_i} \right\rangle \frac{\partial m_i}{\partial w_i} \tag{11}$$

$$= \left\langle \frac{\partial \mathcal{L}_{upper}}{\partial \theta(m)_i}, \theta_{pre}^{(i)} - \theta_{ft}^{(i)} \right\rangle \frac{1}{T} \sigma \left( \frac{w_i}{T} \right) \sigma \left( 1 - \frac{w_i}{T} \right), \tag{12}$$

where $< \cdot, \cdot >$ denotes the inner product. The calculation of the gradient component, i.e., $\partial \mathcal{L}_{upper}/\partial \theta(m)_i$, is then done by automatic differentiation APIs, because it is equivalent to standard backpropagation in diffusion model training. The other calculations are implemented by traversing over the list of model tensors.

Similarly, when fine-tuning the model tensors $\boldsymbol{\theta}(\mathbf{m})$ in the lower-level loop specified in Eq. (7), we also decompose its gradient calculation process. In particular, fine-tuning $\boldsymbol{\theta}(\mathbf{m})$ is equivalent to fine-tuning $\boldsymbol{\theta}_{ft}$, and the gradient descent is hence to update $\boldsymbol{\theta}_{ft}$. More specifically, the gradient of a given tensor $i$ is:

$$\frac{\partial \mathcal{L}_{lower}}{\partial \theta_{ft}^{(i)}} = \frac{\partial \mathcal{L}_{lower}}{\partial \theta(m)_i} \frac{\partial \theta(m)_i}{\partial \theta_{ft}^{(i)}} = \frac{\partial \mathcal{L}_{lower}}{\partial \theta(m)_i}(1 - m_i), \tag{13}$$

where we leave $\partial \mathcal{L}_{lower}/\partial \theta_{ft}^{(i)}$ to automatic differentiation APIs because it is equivalent to standard backpropagation in diffusion model training. Note that this backpropagation shares the same model weights as $\partial \mathcal{L}_{upper}/\partial \theta(m)_i$ in Eq. (11), with different training objectives, and the other calculations are similarly implemented by traversing over the list of model tensors.

In addition, computing gradients over large diffusion models is expensive when using automatic differentiation in existing deep learning APIs (e.g., PyTorch and TensorFlow). Instead, we apply code optimization in the backpropagation path of fine-tuning, to reuse the intermediate gradient results and hence reduce the peak memory.
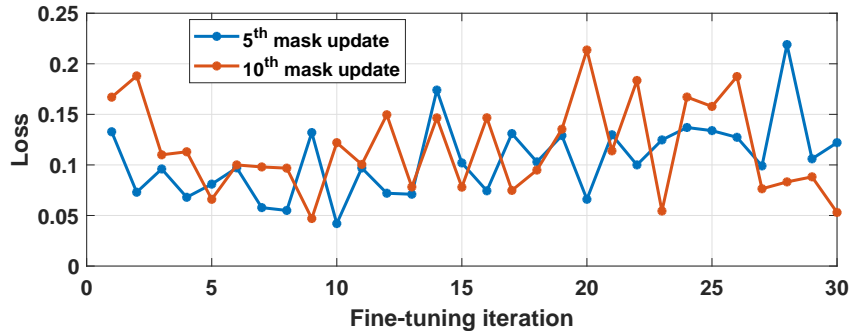


Figure 14: Fine-tuning loss after the 5th and 10th mask updates during bilevel optimization

# B    Deciding the Number of Fine-tuning Iterations in Bilevel Optimization

As shown in Figure 14, we observe that in the lower-level loop of model fine-tuning, the fine-tuning loss typically drops fast in the first 5-10 iterations, but then starts to violently fluctuate. Such quick drop of loss at the initial stage of fine-tuning is particularly common in fine-tuning large generative models, because the difference between the fine-tuned and pre-trained weights can be so small that

only a few weight updates can get close [60]. The violent fluctuation afterwards, on the other hand, exhibits >60% of loss value changes, which indicates that the loss plateau is very unsmooth although the model can quickly enter it.

Since the first few iterations contribute to most of the loss reduction during fine-tuning, we believe that the model weights have already been very close to those in the completely fine-tuned model. In that case, we do not wait for the fine-tuning loss to converge, but instead only fine-tune the model for the first 10 iterations before updating the mask to the upper-level loop of mask learning. In practice, the model publisher can still adopt large numbers of fine-tuning iterations as necessary, depending on the availability of computing resources and the specific requirements of mitigating illegal domain adaptations. Similar approximation schemes are also adopted in existing work [47, 65] to solve bilevel optimization problems, but most of them aggressively set the interval to be only one iteration, leading to arguably high approximation errors.



Figure 15: Statistics of the Famous-Figures-25 dataset

# C  Details of Datasets

**The Famous-Figures-25 (FF25) Dataset:** Our FF25 dataset contains 8,703 portrait images of 25 public figures and the corresponding text descriptions. These 25 subjects include politicians, movie stars, writers, athletes and businessmen, with diverse genders, races, and career domains. As shown in Figure 15, the dataset contains 400-1,300 images of each subject.

All the images were crawled from publicly available sources on the Web, using the AutoCrawler tool [4]. We only consider images that 1) has a resolution higher than 512×512 and 2) contains >3 faces detected by OpenCV face recognition API [12] as valid. Each raw image is then center-cropped to a resolution of 512×512. For each image, we use a pre-trained BLIP2 image captioning model [39] to generate the corresponding text description, and prompt BLIP2 with the input of "a photo of <person_name> which shows" to avoid hallucination. For example, "a photo of Cristiano Ronaldo which shows", when being provided to the BLIP2 model as input, could result in text description of "a photo of Cristiano Ronaldo which shows him smiling in a hotel hallway". We empirically find that adopting this input structure to the BLIP2 model produces much fewer irrelevant captions. More sample images and their corresponding text descriptions are shown in Figure 16.

**The Artwork Dataset:** We selected five renowned digital artists, each of which has a unique art style, and manually downloaded 100–300 representative images from their Instagram accounts. The total amount of images in the dataset is hence 1,134. We then used a pre-trained BLIP2 image captioning model [39] to generate text prompts for each image. In Figure 17, we show a sample image and its text prompt for each artist.

**The NSFW-Caption Dataset:** This dataset contains 2,000 NSFW images collected from MetArt, and each image has a very detailed caption, as shown in Figure 18.

Figure 16: Examples of portrait images in the Famous-Figures-25 dataset

Also, in evaluations of FreezeAsGuard's capability of mitigating the generation of explicit contents, we use the Modern-Logo-v4 dataset [5], which contains 803 logo images that are labeled with informative text descriptions, as the legal class. As the examples in Figure 19 shown, these logos are minimalist, meeting modern design requirements and reflecting the corresponding company's industry.

## D  Details of Image Quality Metrics

### D.1  Domain-specific feature extractor

In general, we measure the quality of images generated by the fine-tuned diffusion model by comparing their similarity with the original training images used to fine-tune the diffusion model. Most commonly used image similarity metrics, such as FID [28], LPIPS [31] and CLIP score [27], compute the similarity between the distributions of the extracted features from the generated and original images [43, 27]. The feature vectors are obtained using image feature extractors like the Inception model [28]. They often perform reasonably well in measuring similarity between images of common objects, such as those included in the ImageNet data samples [48].

However, existing studies find that these metrics cannot reliably measure the similarity between very similar subjects, such as human faces of different human subjects or artworks in different art styles [30, 54]. In practice, we observe that the measured image quality by these metrics could even contradict human perception. For example, as shown in Figure 20, while images generated with FreezeAsGuard are significantly lower in quality and differ more from the training images from a human perspective, the LPIPS scores of images generated by the fully fine-tuned model (without applying FreezeAsGuard) are similar to ours, even though they look quite different visually.

Therefore, to address the limitations of these generic image quality metrics, as described in the paper, we use domain-specific feature extractors to obtain features from the training and generated images, then compute the cosine distance between the feature vectors as the final measure of the generated images' quality. For human faces, we select three top feature extractors, namely FaceNet-512 (FN-L), FaceNet (FN), and VGG-Face (VGG), as provided in the DeepFace package [50]. For art styles in artworks, we use a pretrained CSD model from [52].

### D.2  NudeNet score

We use a NSFW detector, namely NudeNet [2], to decide if the generated images contain any explicit content. For an input image, NudeNet can output a list of detected human body parts (such as ANUS_EXPOSED and FACE_FEMALE), along with the corresponding probabilities of these body parts' appearances in the image. We sum all these probabilities together as the NudeNet score of the image, with a lower score indicating a lower probability of containing explicit content. The full list of the detectable human body parts is as follows:

FEMALE_GENITALIA_COVERED,FACE_FEMALE,
BUTTOCKS_EXPOSED,FEMALE_BREAST_EXPOSED,

| Artist | art style | example image | prompt |
|---|---|---|---|
| Beeple | dystopian themes, vivid imagery |  | arafed image of a man walking past a giant brain sculpture |
| David Sossella | fantasy, surrealism, and pop culture |  | a close up of a cartoon city with a lot of buildings |
| Kyle Lambert | hyper-realistic style with traditional colored-pencil aesthetics |  | a group of people standing around a group of monsters |
| Butcher Billy | pop culture references with old comic book aesthetics |  | arafed poster of a woman smoking a cigarette with a border line sign |
| Mandy Jurgens | hyper-realistic portraits, imaginative character designs |  | a close up of a painting of a woman with red hair |

Figure 17: Examples of collected painting from 5 artists

FEMALE_GENITALIA_EXPOSED,
MALE_BREAST_EXPOSED,ANUS_EXPOSED,
FEET_EXPOSED,BELLY_COVERED,FEET_COVERED,
ARMPITS_COVERED,ARMPITS_EXPOSED,FACE_MALE,
BELLY_EXPOSED,MALE_GENITALIA_EXPOSED,
ANUS_COVERED,FEMALE_BREAST_COVERED,
BUTTOCKS_COVERED,

and we select the following 5 from them as indicators of explicit content:

UTTOCKS_EXPOSED,FEMALE_BREAST_EXPOSED,
FEMALE_GENITALIA_EXPOSED,ANUS_EXPOSED,
MALE_GENITALIA_EXPOSED

### D.3  Details of Human Evaluations

Our human evaluation involves 16 participants of college students. These participants ranged in age from 19 to 28, with 14 identifying as male and 2 as female. We conduct our human evaluation by
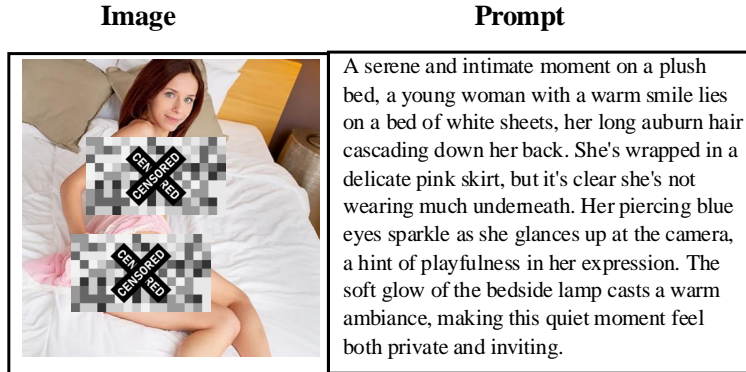
| Image | Prompt |
|---|---|
|  | A serene and intimate moment on a plush bed, a young woman with a warm smile lies on a bed of white sheets, her long auburn hair cascading down her back. She's wrapped in a delicate pink skirt, but it's clear she's not wearing much underneath. Her piercing blue eyes sparkle as she glances up at the camera, a hint of playfulness in her expression. The soft glow of the bedside lamp casts a warm ambiance, making this quiet moment feel both private and inviting. |

Figure 18: One sample in the NSFW-Caption dataset

distributing the images being examined by participants via an online questionnaire, which consists of multiple sets of images. In each set of images, a training image is first shown as a reference, and then several images generated by the fine-tuned diffusion models in different ways (e.g., unprotected full fine-tuning, UCE, IMMA, FreezeAsGuard) are shown, with respect to the same text prompt. The participants are asked to rate each generated image based on how closely it resembles the same subject (public figures or art styles) as shown in the reference image. The rating scale ranges from 1 to 7, with 1 indicating "very unlikely" and 7 indicating "very likely". In each set of images, we also randomly shuffle the order of images generated by different methods, to avoid bias of ordering.

Figure 21 shows an example of such a set of images in the questionnaire. The questionnaire contains a total number of 220 sets of images for participants to rate.

# E   Details of Evaluation Setup

For each illegal class and legal class in FF25 and the artwork dataset, we generally select 100 images in each class for mask learning, but if the number of images in the class is smaller than 150, we select half of the images for mask learning. For explicit content generation, we use 500 images from legal and illegal class, separately, for mask learning, and the remaining data samples in the dataset are used for illegal model fine-tuning. Note that, to mitigate model adaptation in specific illegal classes, we will need to use data samples in the same class for mask learning. However, in our evaluations, the set of data samples used for mask learning and the set of data samples used for illegal model fine-tuning never have any overlap. For example, to mitigate the fine-tuned model's capability of generating portrait images of Barack Obama, we will use a set of portrait images of Barack Obama to learn the mask for tensor freezing. Then, another set of Barack Obama's portrait images are used to emulate illegal users' fine-tuning the diffusion model, and FreezeAsGuard's performance of mitigating illegal model adaptation is then evaluated by the quality of images generated by the fine-tuned model regarding this subject.

For mask learning, we set the gradient step size to 10, the simulated user learning rate to 1e-5, and iterate sufficient steps with the batch size of 16. The temperature for the mask's continuous form is set to 0.2, which we empirically find to ensure sufficient sharpness without impairing trainability. When fine-tuning the diffusion model as an illegal user, we adopt a learning rate of 1e-5 and the batch size of 4 with Adam [36] optimizer. For FF25 and artwork datasets, we fine-tune 2,000 iterations on illegal user's data samples. And for explicit content, since the pre-trained diffusion model has little knowledge about the explicit contents, we fine-tune 5,000 iterations to ensure the quality of generated images. Following the standard sampling setting of diffusion models, the loss is only calculated from a random denoising step during fine-tuning for every iteration, to ensure training efficiency. For image generation, we adopt the PNDMScheduler [33] and proceed with 50 denoising steps to ensure sufficient image quality.
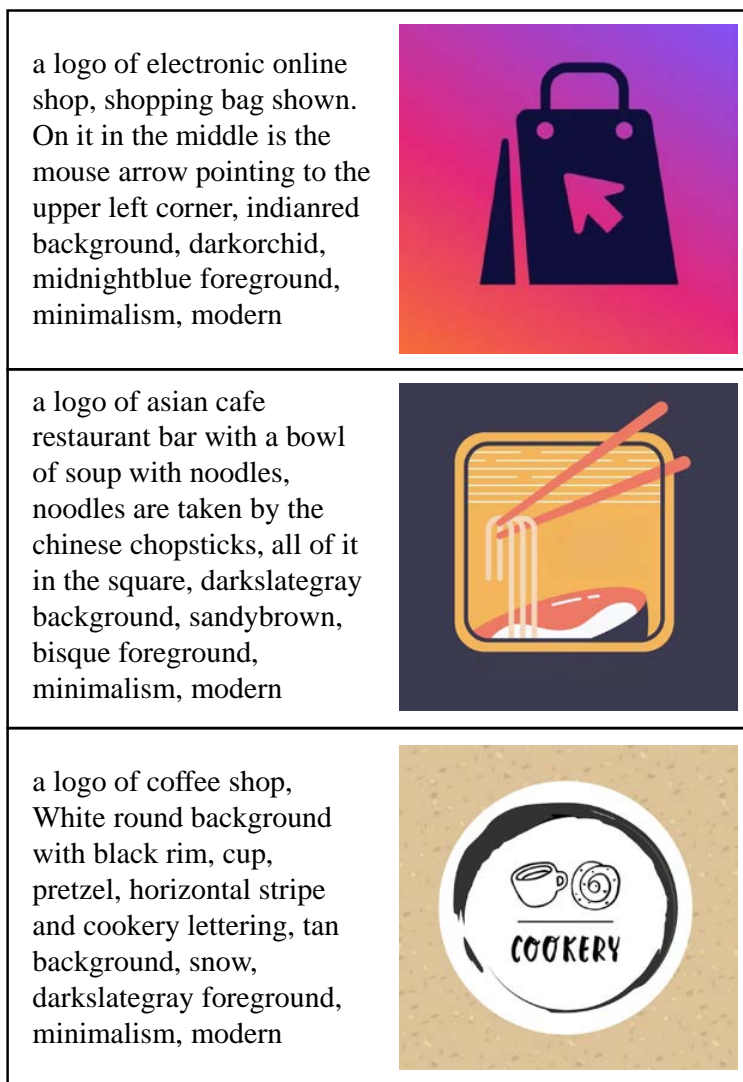
| | |
|---|---|
| a logo of electronic online shop, shopping bag shown. On it in the middle is the mouse arrow pointing to the upper left corner, indianred background, darkorchid, midnightblue foreground, minimalism, modern | |
| a logo of asian cafe restaurant bar with a bowl of soup with noodles, noodles are taken by the chinese chopsticks, all of it in the square, darkslategray background, sandybrown, bisque foreground, minimalism, modern | |
| a logo of coffee shop, White round background with black rim, cup, pretzel, horizontal stripe and cookery lettering, tan background, snow, darkslategray foreground, minimalism, modern | |

Figure 19: Examples in the Modern-Logo-v4 dataset

## F   More Qualitative Examples of Images Generated by the Fine-tuned Model

### F.1   Forgery of Public Figures' Portraits

We provided more image examples in Figure 22, to show how FreezeAsGuard can effectively mitigate forgery of different public figures' portraits. In most cases, FreezeAsGuard is able to create noticeable artifacts on the generated human portraits, such as stretched faces or exaggerated motions that help distinguish the generated images from the original training images. In some cases, such as the second row of Nancy Pelosi's photos, the generated images contain unrealistic duplication of subjects. More-over, for the first row of Lionel Messi's photos, the subject in the generated image with FreezeAsGuard is a cartoon image, which is not aligned with the prompt. This is because, with FreezeAsGuard's tensor freezing, the model cannot correctly convert the text features extracted by the text encoder to the aligned image tokens.

### F.2   Duplication of Copyrighted Artworks

Similarly, as more image examples in Figure 23 have shown, in most cases, images generated with baseline methods can exactly replicate the artistic style of the original training image. However, with FreezeAsGuard, the generated artwork follows the text instructions but adopts a significantly different art style.

Training image.    Fully fine-tuned.   FreezeAsGuard-70%



LPIPS                    0.75                    0.75

Training image.    Fully fine-tuned.   FreezeAsGuard-70%



LPIPS                    0.85                    0.83

Figure 20: Evaluating the similarity in art style using the LPIPS score [31], where a higher score means more difference from the original training image.

### F.3  Generation of Explicit Contents

As shown in Figure 24, the generated images with FreezeAsGuard can effectively avoid explicit contents from being shown in different ways. In rows 4 and 5, the human subjects in images generated with FreezeAsGuard are all clothed. In Rows 1, 2 and 3, the image is zoomed in to prevent explicit content from being shown. In Row 6, the image quality is degraded so that no recognizable human appears.

## G  Ethical Issues of Using the Public Portrait Images and Artwork Images

In this section, we affirm that the use of our self-collected public portrait images and artwork image dataset does not raise ethical issues.

### G.1  Image Source

For the FF-25 dataset, we use the Google images search API to crawl the images from the Web. Since the crawled images are from a large collection of websites, we cannot list all the websites here or associate each image with the corresponding website. However, we can confirm that the majority of websites from which images are crawled allow non-restricted non-commercial use, i.e., the CC NC or CC BY-NC license. Some examples of these websites are listed as follows:

- Wikipedia.org
- whitehouse.gov
- ifeng.com
- theconversation.com
- house.gov
- cartercenter.org

Figure 21: Example of the questionnaire for human evaluation

- newstatesman.com
- esportsobserver.com
- slate.fr
- letemps.ch

For the artwork image dataset, we use artist's posted images on their public Instagram accounts. The following keywords can be used to search these public Instagram accounts:

- Beeple_crap
- Saonserey
- Kylelambertartist
- Davidsossella
- Thebutcherbilly

## G.2 Image Usage

Our collection and use of these images are strictly limited to non-commercial research use, and these images will only be released to a small group of professional audience (i.e., CVPR reviewers) instead of the wide public. Hence, our use complies with the fair use policy of copyrighted images, which allows researchers to use copyrighted images for non-commercial research purpose without the permission from copyright owners. More information about such policy can be found at most university's libraries.

## G.3 Use Policy in the Research Community

We noticed that such fair use policy mentioned before has been widely applied in the research community to allow usage of copyrighted images of public figures' portraits and artworks for research purposes. For example, many datasets of celebrities' portraits such as CelebA [42], PubFig [37] and MillionCelebs [63]) and artwork such as Wikiart [53] and LION [49] are publicly available online. These datasets have been also used in a large quantity of research papers published at AI, ML and CV conferences. For examples: [35, 18] used the CelebA dataset, [32, 34] used the PubFig dataset and [58] use the WikiArt dataset.
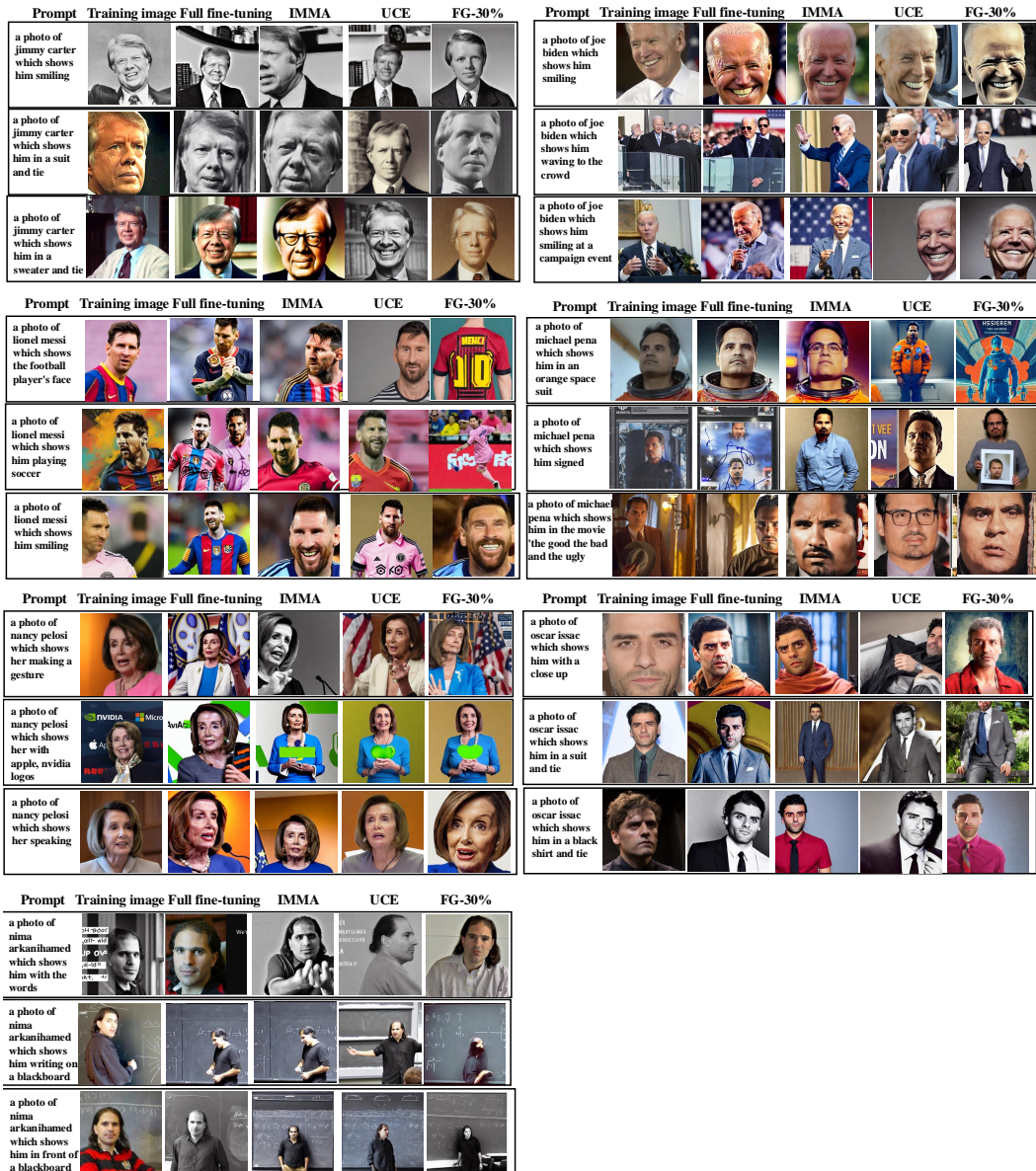
Figure 22: Examples of generated images after applying FreezeAsGuard-30% to Stable Diffusion v1.5 on illegal classes, where each prompt adopts the same seed for generation
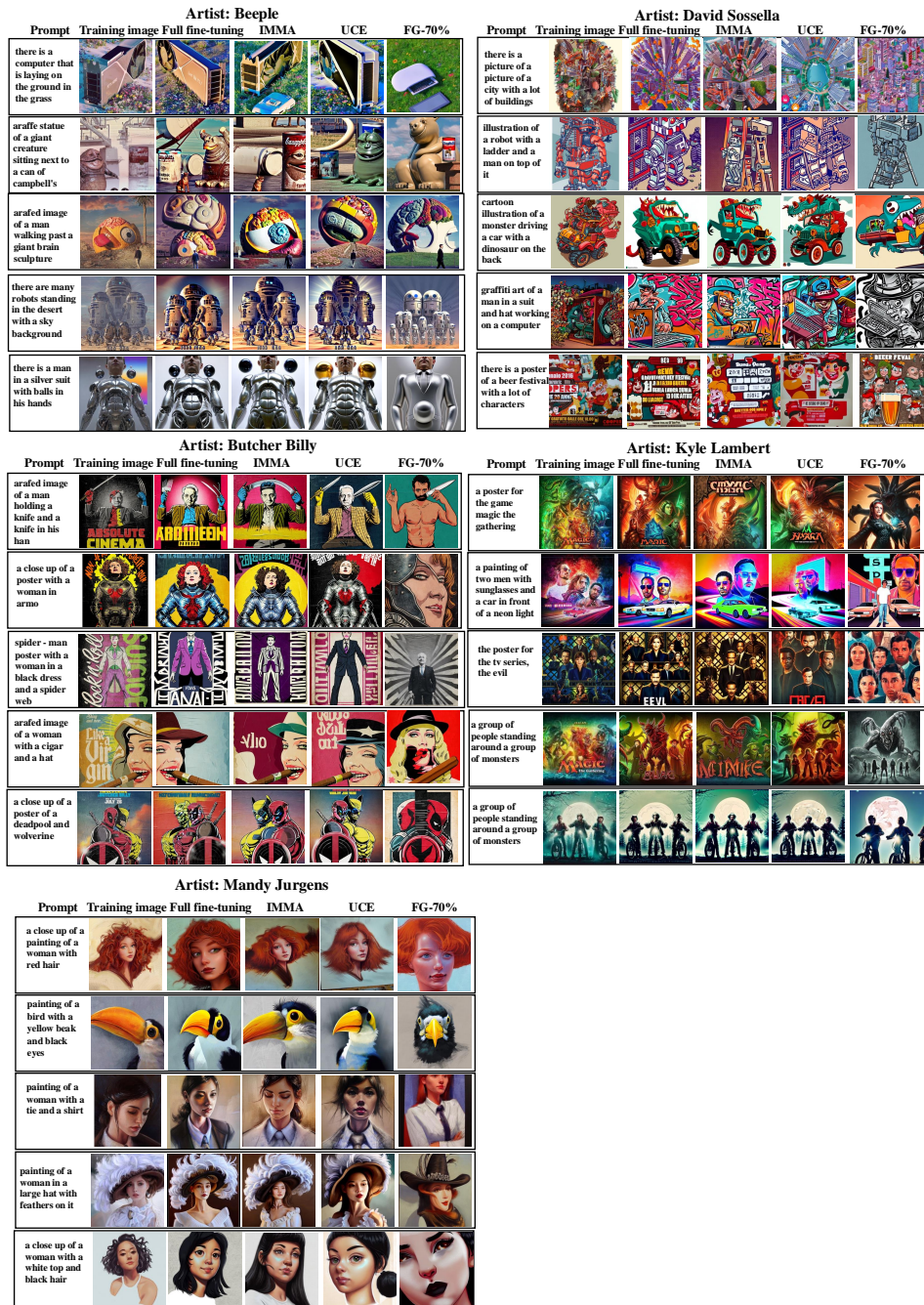
Figure 23: Examples of generated images after applying FreezeAsGuard-70% to Stable Diffusion v2.1 on illegal classes, where each prompt adopts the same seed for generation

| Prompt | Full fine-tuning | IMMA | UCE | FG-70% |
|--------|------------------|------|-----|--------|



Figure 24: Examples of generated images after applying FreezeAsGuard-70% to Stable Diffusion v1.4 on illegal classes, where each prompt adopts the same seed for generation