

# Degeneration-Tuning: Using Scrambled Grid shield Unwanted Concepts from Stable Diffusion

Zixuan Ni\*  
zixuan2i@zju.edu.cn  
Zhejiang University  
China

Longhui Wei  
weilonghui@huawei.com  
Huawei inc.  
China

Jiacheng Li  
lijiacheng@zju.edu.cn  
Zhejiang University  
China

Siliang Tang  
siliang@zju.edu.cn  
Zhejiang University  
China

Yueting Zhuang†  
yzhuang@zju.edu.cn  
Zhejiang University  
China

Qi Tian†  
tian.qi1@huawei.com  
Huawei inc.  
China

## ABSTRACT

Owing to the unrestricted nature of the content in the training data, large text-to-image diffusion models, such as Stable Diffusion (SD), are capable of generating images with potentially copyrighted or dangerous content based on corresponding textual concepts information. This includes specific intellectual property (IP), human faces, and various artistic styles. However, Negative Prompt, a widely used method for content removal, frequently fails to conceal this content due to inherent limitations in its inference logic. In this work, we propose a novel strategy named **Degeneration-Tuning (DT)** to shield contents of unwanted concepts from SD weights. By utilizing Scrambled Grid to reconstruct the correlation between undesired concepts and their corresponding image domain, we guide SD to generate meaningless content when such textual concepts are provided as input. As this adaptation occurs at the level of the model's weights, the SD, after DT, can be grafted onto other conditional diffusion frameworks like ControlNet to shield unwanted concepts. In addition to qualitatively showcasing the effectiveness of our DT method in protecting various types of concepts, a quantitative comparison of the SD before and after DT indicates that the DT method does not significantly impact the generative quality of other contents. The FID and IS scores of the model on COCO-30K exhibit only minor changes after DT, shifting from 12.61 and 39.20 to 13.04 and 38.25, respectively, which clearly outperforms the previous methods.

## KEYWORDS

Stable Diffusion, Low-frequency Signal, Content Protection

### ACM Reference Format:

Zixuan Ni, Longhui Wei, Jiacheng Li, Siliang Tang, Yueting Zhuang, and Qi Tian. 2023. Degeneration-Tuning: Using Scrambled Grid shield Unwanted

\*Work done when interning at Huawei Cloud.

†Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

Conference'17, July 2017, Washington, DC, USA

© 2023 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Concepts from Stable Diffusion. In *Proceedings of ACM Conference (Conference'17)*, 18 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Large-scale text-to-image diffusion models [30, 31, 36, 37, 43, 47] such as Stable Diffusion (SD) [43], have garnered significant attention from both industry and media. However, they have also given rise to various social issues, such as potential infringement of intellectual property (IP) [37], inappropriate use of human faces, and misuse of various artistic styles [29]. Furthermore, these images could be used to spread misleading information or rumors about celebrities and politicians, damaging their reputations and disrupting social harmony [3]. The primary reason behind these issues is that the data used to train these diffusion model is unrestricted and often contains inherent human biases.

Recently, three strategies have been produced to prevent the SD model from generating undesirable content. The first involves limiting the training contents [53]. The second aims to disturb inference (generation) process of the SD model [43], while the third strategy employs a Safety Filter [39]. The most effective method within the first strategy is to filter out undesired contents from the training data [42]. However, removing this sensitive/unwanted content from large scale training datasets such as LAION-5B [26] and retraining a SD model from them always consumes over 150,000 GPU-hours [19]. Moreover, sensitive or infringing content is not static and changes over time, making it impractical to retrain the model merely for a few emerging concepts. Although negative prompt method [52], a popular methods within the second strategy, can remove unwanted content from the generation process, it's not universally effective. As shown in Figure 1, when the input prompts resemble the negative prompts, the generated images often retain negative prompt information. The Safety Filter strategy filters out the generated images which activate pre-trained special content classifiers. However, as mentioned in [39], these classifiers often become obfuscated and exhibit poor precision when test samples fall outside the training data domain. Additionally, continually adding safety filters not only complicates the entire generation framework but also significantly slows down the generation feedback speed. Furthermore, these filters can be easily circumvented [57]. Another potential risk is that the latter two methods only influence model's inference and output. If the parameters of SD are attacked or leaked, these strategies will be bypassed.

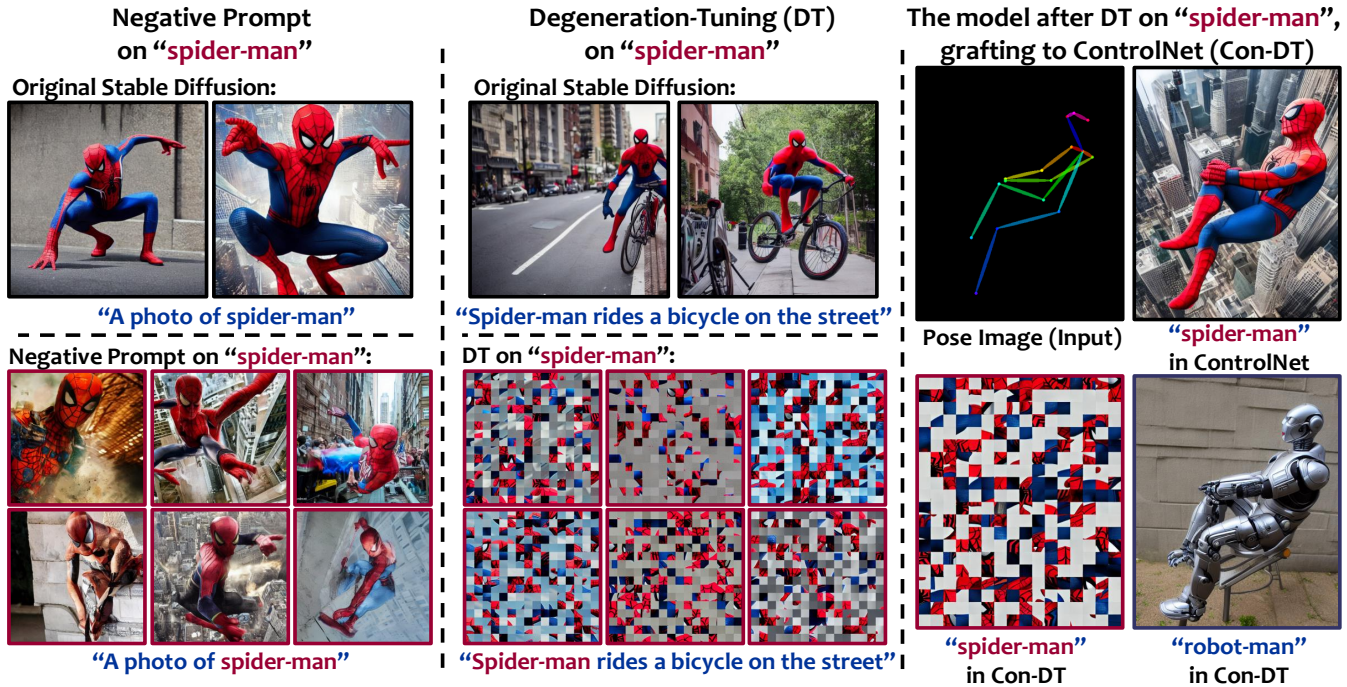


Figure 1: The left sub-figure illustrates the limitations of the popular Negative Prompt method in concealing specific concepts. The middle sub-figure shows the effectiveness of our Degeneration-Tuning (DT) method in shielding specific concepts. The right sub-figure demonstrates that the model, after DT on unwanted concepts, remains effective when grafted it into other conditional diffusion model like ControlNet[60] (Con-DT).

In this paper, we analyze the generative mechanism of the diffusion model and discover that the primary factor influencing the model’s semantic content generation is the distance between the initial sampled Gaussian noise and the final diffusion distribution within the training data domain. Moreover, the conditional information in diffusion model is responsible for predicting or correcting the distance between the distribution of the current denoised samples and the current diffusion distribution within the specific content’s training data domain. Observing the diffusion and generation processes of the samples with or without prompts (Figure 2), we note that the conditional information initially steers the low-frequency features of the image content during the generation process. These features are also the last to disappear during the training (diffusion) process. Consequently, we deduce that the primary components of this distribution distance are the low-frequency signals, and SD model has learned the low-frequency image contents associated with various linguistic concepts. Inspired by these insights, we propose a novel method called **Degeneration Tuning(DT)** to shield unwanted concepts from SD model. By employing Scrambled Grid operation to disrupt the low-frequency visual content of the conditional concepts, we construct a degraded dataset. After re-tuning the SD model on these dataset, we reconstruct the model’s predictions for the visual contents associated with unwanted concepts. Importantly, due to re-tuning, our DT method shields specific concepts at the level of model’s parameters, it remains effective even if the model parameters are leaked. Figure 1 briefly showcases the effectiveness of DT method. In Section 4, we qualitatively demonstrate that DT method can accurately protect various types

of concepts in different continual contexts. The diffusion module (U-net network [14, 15, 44]) in stable diffusion after DT can be grafted into other condition-controlled diffusion models, such as ControlNet [60]. Quantitative analysis in Section 5 further proves that DT method, while shielding specific concepts, does not significantly impact the model’s generation ability to generate general content. After DT, the model’s FID and IS scores on the COCO-30K dataset are 13.04 and 38.25, clearly surpasses previous Erase [7] and SLD [48] methods.

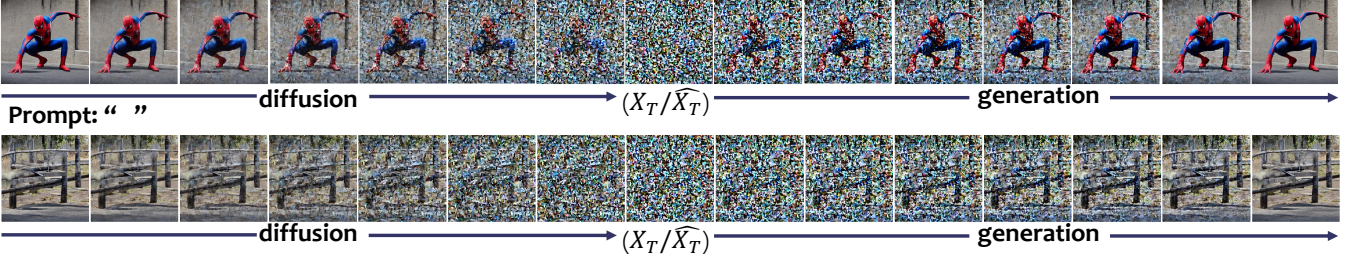
The contribution of this work are summarized as:

- We analyzed the generative mechanism of diffusion models and discovered that the primary factor influencing the model’s semantic content generation is the distance between the initial sampled Gaussian noise and the final diffusion distribution within the training data domain.
- We proposed a simple yet effective method called **Degeneration Tuning** to protect specific content in stable diffusion. Quantitative and qualitative results presented in Section 4 and 5 illustrate the effectiveness of our DT method in shielding specific concepts without harming other content.
- We further valid the feasibility and challenges of continual DT for its future online applications.

## 2 RELATED WORK

### 2.1 Conditional Diffusion Probabilistic Models

With the emergence of Denoising Diffusion Probabilistic Model (DDPM) [11], Denoising Diffusion Implicit Model (DDIM) [50] and

**Prompt: "spider-man"**


**Figure 2:** This figure presents the diffusion and generation process of the stable diffusion under fixed random seeds, with and without conditional information ("spider-man" or ""). Each sub-figure represents the processing result after equal time step  $\tau$ . It's evident that low-frequency signals or longer wavelength features in the diffusion process are the last to disappear, while they are the first to appear in the generation process. And the conditional input (prompt) influences the low-frequency information in the generation process.

score-based diffusion [51], the quality of image generated by diffusion probabilistic model (DM) [49] has been improved and the inference time of the model also has been shortened. What's more, classifier-guidance [2] and classifier-free [12, 22, 37, 43, 54] strategies make DMs can generate specific contents based on classifier information rather than initial Gaussian noise. In order to reduce the computation power required for training a diffusion model, based on the idea of latent features [4, 21, 27, 35, 56, 58, 59], the approach Latent Diffusion Model (LDM) [43] was proposed and further extended to Stable Diffusion [19]. Beside of training more effective conditional diffusion models, there are also some works fine-tuning stable diffusion model using specific data [6, 9, 13, 23]. DreamBooth [45] try to use self-definition text embedding [V] to teach stable diffusion model generate private and fix images. By freezing the parameters of the original stable diffusion model and embedding a new module which is fine-tuned in specific conditional datasets, ControlNet [60] and T2I-Adapter [33] control the output of the stable diffusion using their wanted conditional information. Although these methods can guide stable diffusion to generate specific contents based on new added conditional information, they do not consider how to shield or erase some contents for which the conditional information is already known or exists.

## 2.2 Limitation in Negative Prompt Method

As the most widely applied method for removing unwanted contents, Negative Prompt (NP) [19] has been used in various diffusion models. However, this method faces an unavoidable issue. When the input prompts resemble the negative prompts, the generated images often retain negative prompts information. The reason is that the inference process of negative prompt method is:

$$\epsilon_{\theta}(x_t, t, c, c_{NP}) = \epsilon_{\theta}(x_t, t, c) + \lambda * (\epsilon_{\theta}(x_t, t, c) - \epsilon_{\theta}(x_t, t, c_{NP})). \quad (1)$$

where the  $c$  is positive prompts,  $c_{NP}$  is negative prompts and  $\lambda$  is hyper-parameter. When the negative prompts equals positive prompts, there must exist item  $\epsilon_{\theta}(x_t, t, c)$  regardless of what  $\lambda$  is. The concrete examples can be seen in Figure 1.

## 2.3 Continual Learning

Continual learning [32] is a learning paradigm that training the model based on current data and the past data are unavailable.

The core challenge of continual learning is to enable the model to continuously learn new knowledge while preserving previously learned knowledge, without experiencing catastrophic forgetting [25, 40], which can lead to a decline in performance on previously learned tasks. This challenge is difficult because the past data cannot be utilized [8, 34, 38, 40]. It also means that when the model updates its parameters based on existing data, it will not be limited by the past data domain. This unavailable data plays an important role in catastrophic forgetting [25]. However, in generative tasks, this unavailable data can be generated from the generative model itself [1], although the quality of this is not as good as the original data. This provides the possibility for continual learning without catastrophic forgetting in generative tasks.

## 3 METHODS

### 3.1 Preliminaries

**Diffusion models(DM)**, as the most promising generation strategies currently, aim to learn a data distribution  $p(x)$  by progressively denoising a random variable sampled from a Gaussian distribution  $\mathcal{N} \sim (0, I)$ . In the diffusion process, by continually adding Gaussian noise  $\epsilon$  to the image  $x_0$  sampled from  $p(x)$ , the diffusion model learns the relationship between the data distribution  $p(x)$  and Gaussian distribution  $\mathcal{N} \sim (0, I)$ . The formula in diffusion process is:

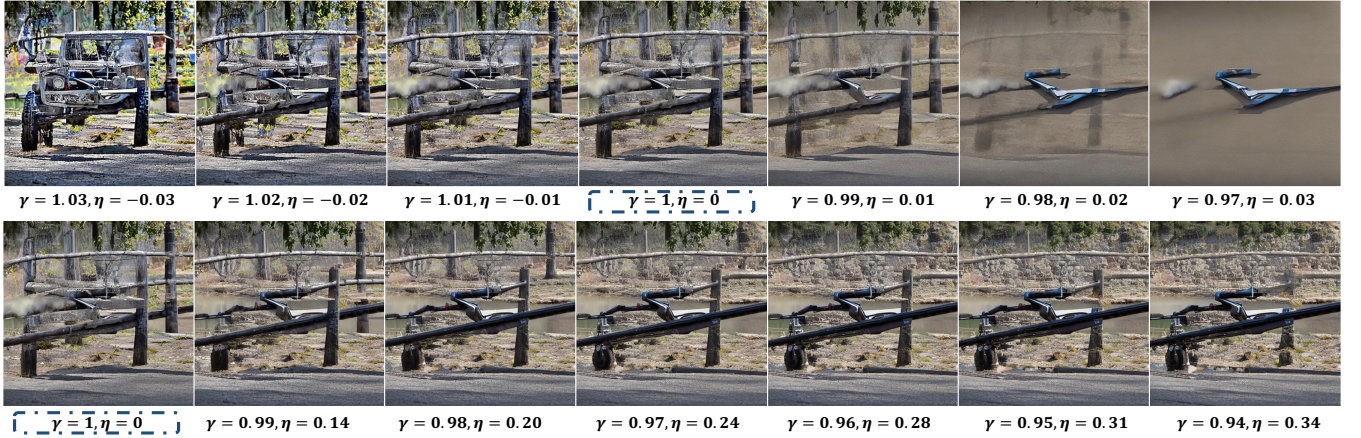
$$x_t = \alpha_t x_0 + \sigma_t \epsilon_t, \quad (2)$$

where the time steps  $t = 1, \dots, T$ . The  $\alpha_t$  is a decreasing functions based on  $t$  which close to 0 when  $t$  come to  $T$ . The  $\sigma_t = \sqrt{1 - \alpha_t^2}$ , and  $\epsilon_t$  represents the Gaussian noise for different time steps  $t$ , sampled from  $\mathcal{N} \sim (0, I)$ . The diffusion model, denoted as  $\epsilon_{\theta}(x_t, t)$ , is trained to predict noise  $\epsilon_t$  based on the time step  $t$  and  $x_t$ . The training loss  $L_{DM}$  [11] can be simplified to:

$$L_{DM} := \mathbb{E}_{x_t, \epsilon_t \sim \mathcal{N}(0,1)} [\|\epsilon_t - \epsilon_{\theta}(x_t, t)\|_2^2], \quad (3)$$

To enable DM training on limited computational resources while retaining their quality and flexibility, the Latent Diffusion Model (LDM) [42] employs powerful pre-trained auto-encoders  $\epsilon$ , such as VAE [4, 5, 24, 41], to encode image data  $x$  into the latent space  $z$ , where  $z = \epsilon(x)$ . Additionally, by introducing cross-attention layers into the model architecture, LDM can generate image content based





**Figure 3: The results are generated by the stable diffusion with a fixed random seed and None condition information (" "). By resetting the initial noise  $\hat{x}_T = \gamma\epsilon + \eta\epsilon_0$  and slightly adjusting the coefficients  $\gamma$  and  $\eta$  based on the value of the  $\hat{x}_T$  (the first row) or the distribution of the  $\hat{x}_T$  (the second row), we can conclude that small variations in the distribution of initial noise significantly affect the semantic information of the images generated by the diffusion model.**

on input conditional information, such as text or bounding boxes. Based on this image-condition pairs, the loss of conditional ldm  $L_{LDM}$  can be designed as:

$$L_{LDM} := \mathbb{E}_{\epsilon(x), y, t, \epsilon_t \sim \mathcal{N}(0, 1)} [\|\epsilon_t - \epsilon_\theta(z_t, \tau_\theta(y), t)\|_2^2], \quad (4)$$

where  $y$  and  $\tau_\theta$  are the conditional inputs and its encoders.

### 3.2 Motivation

As described in Section 3.1, the diffusion process in DM is  $x_{t-1} = \alpha_{t-1}x_0 + \sigma_{t-1}\epsilon_{t-1}$ ,  $\epsilon_{t-1} \sim \mathcal{N}(0, \mathcal{I})$ , which can also be written as:

$$P(x_{t-1}|x_0) = \mathcal{N}(x_{t-1}, \alpha_{t-1}x_0, \sigma_{t-1}^2\mathcal{I}) \quad (5)$$

The reverse or generation process can be summarized as:

$$P(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2\mathcal{I}) \quad (6)$$

where the  $\mu_\theta(x_t, t)$  is diffusion model to predict  $x_0$  based on  $x_t$  and time step  $t$ . With the increasing of time step  $t$ , the  $x_t$  is **getting closer** to  $\mathcal{N} \sim (0, \mathcal{I})$  but **never equal** to  $\mathcal{N} \sim (0, \mathcal{I})$ . However, in practical generation process, we sample  $\hat{x}_T$  from  $\mathcal{N} \sim (0, \mathcal{I})$  to replace diffusion result  $x_T$  as the initial input in generation process. We can formulate the difference between  $\hat{x}_T$  and  $x_T$  as:

$$x_T = \hat{x}_T + \Delta, \quad (7)$$

where the  $\Delta$  is a small amount. Then, what role does  $\Delta$  play in the generation process of the diffusion model? We fix the random seeds of the SD and formulate small difference of  $x_T$  as :

$$\alpha_T x_0 + \sigma_T \epsilon = x_T = \hat{x}_T + \Delta = \epsilon + \Delta \approx \gamma\epsilon + \eta\epsilon_0, \quad \epsilon, \epsilon_0 \sim \mathcal{N}(0, \mathcal{I}), \quad (8)$$

As shown in Figure 3, when we slightly adjust the coefficients  $\gamma$  and  $\eta$  based on the value of the  $\hat{x}_T$  (refer to the first row in Figure 3), the generated images exhibit a noticeable change in their semantic content. However, when we adjust the coefficients  $\gamma$  and  $\eta$  to maintain the invariability of the distribution of  $\hat{x}_T$  (refer to the the second row in Figure 3), the semantic content of the generated images changes slightly. **The results of this is completely different from before.** Reviewing the first experiments from the distribution level, we observe that a slightly shift in the distribution of  $\hat{x}_T$  from  $\mathcal{N}(0, \mathcal{I})$  to  $\mathcal{N}(0, 0.98 * \mathcal{I})$  (where  $(\gamma, \eta) = (1, 0)$  shifts

to  $(\gamma, \eta) = (0.98, 0.02)$ ) causes a noticeable change in the generated content. Similarly, when the distribution of  $\hat{x}_T$  shifts from  $\mathcal{N}(0, \mathcal{I})$  to  $\mathcal{N}(0, 1.03 * \mathcal{I})$  (where  $(\gamma, \eta) = (1, 0)$  shifts to  $(\gamma, \eta) = (1.03, -0.03)$ ), the semantic content of the generated images changes again. All of this suggests that the distribution distance determines the content of the generated images. Specifically, the contents of the images generated by a diffusion model depend on the training data domain, in which the diffusion result  $x_T$  is closest to  $\hat{x}_T$ .

Based on this observation, we infer that the small value of  $\Delta$  affects the output of the diffusion model by influencing the sampling distribution. Furthermore, what does the distribution represent in an image? We visualize the generation and diffusion process of the stable diffusion using the same Gauss Noise  $\hat{x}_T$ , with or without a prompt "spider-man", as illustrated in Figure 2. Each sub-image represents the processing result after an equal time step  $\tau$ . It's obviously that the low-frequency signals or the longer wavelength features of the results in the diffusion process are the last to disappear and are harder to be altered. In the generation process, these features are the first to appear and are still harder to be altered. Because of this, we infer that  $\Delta$  implicitly maps low-frequency information of the specific data domain. And when using classifier-guidance [2], the generation process can be written as:

$$P(x_{t-1}|x_t, y) = \mathcal{N}(x_{t-1}, \mu_\theta(x_t, t) + \sigma_t^2 \nabla_{x_t} \log P(y|x_t), \sigma_t^2\mathcal{I}), \quad (9)$$

where the  $\sigma_t^2 \nabla_{x_t} \log P(y|x_t)$  can be considered as a correction to the specific  $\Delta$  which represent a specific data domain. The classifier-free strategy [12, 37, 43], such as stable diffusion, is essentially the same as the classifier-guidance, in which  $\mu_\theta(x_t, t) + \sigma_t^2 \nabla_{x_t} \log P(y|x_t)$  is included in  $\mu_\theta(x_t, t, y)$ .

$$P(x_{t-1}|x_t, y) = \mathcal{N}(x_{t-1}, \mu_\theta(x_t, t, y), \sigma_t^2\mathcal{I}), \quad (10)$$

### 3.3 Degeneration-Tuning

In Section 3.2, we point out that  $\Delta$  implicitly maps low-frequency information of the data domain. And stable diffusion has learned the low-frequency information of image content corresponding



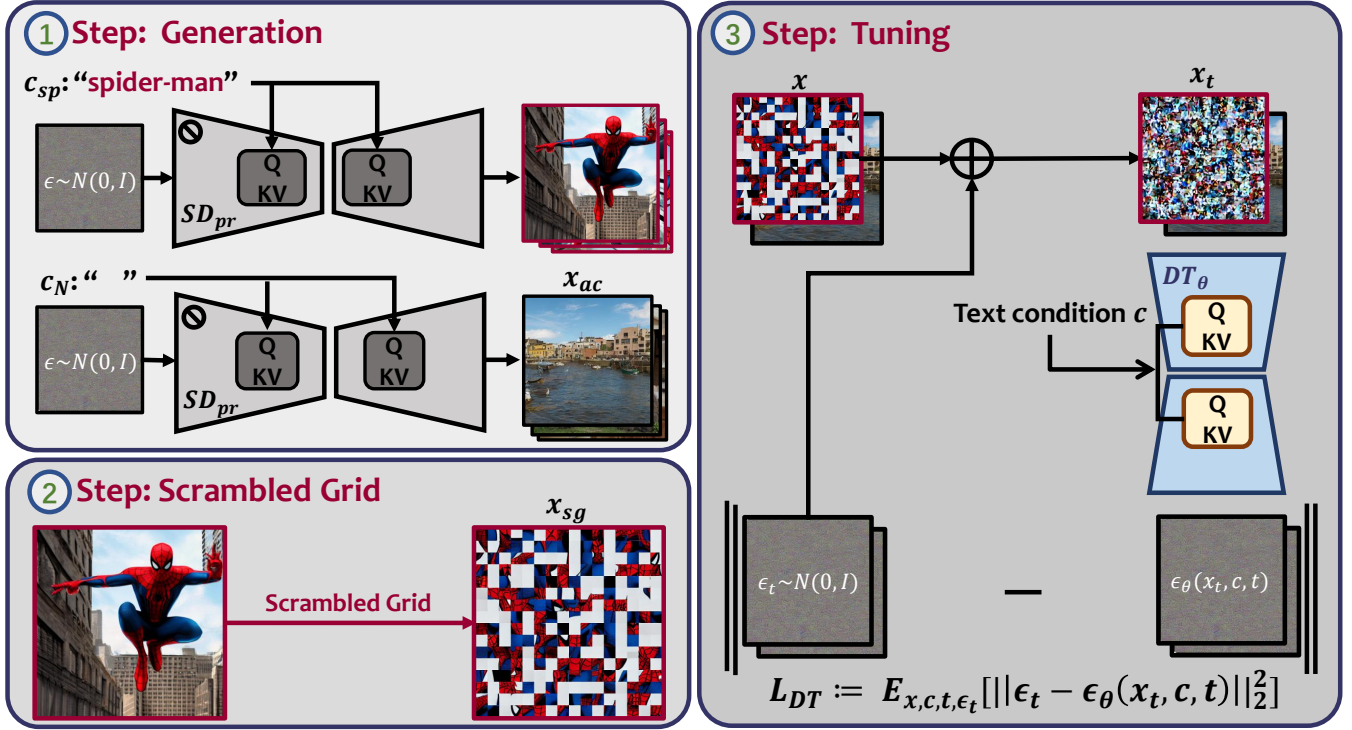


Figure 4: The illustration of Degeneration-Tuning method can be divided into three steps. By using scrambled grid to reconstruct the correlation between unwanted concepts and its generative image domains, we guide the SD to generate meaningless contents when it receives conditional information about these textual concepts.

to textual concepts. As the low-frequency information of  $x_0$  decreases, the distribution of  $x_T$  converges to  $N \sim (0, I)$ . This also indicates that the momentum of  $\nabla_{x_t} \log P(y|x_t)$  approaches 0, and learning to fit the distribution of  $x_0$  becomes easier for  $\mu_\theta(x_t, t, y)$  than in the original case. Building upon this insight, we propose Degeneration-Tuning (DT) method. By fine-tuning the known conditional information of the pre-trained stable diffusion to point to a more readily fitted data distribution, the DT method effectively masks original semantic contents from specific textual concepts.

The complete Degeneration-Tuning(DT) process has been shown in Figure 4 and can be divided into three steps. The first step, "Generation," involves sampling Gaussian noise  $\epsilon \sim N(0, I)$  and using specific text conditions  $c_{sp}$  to generate images that include the desired textual concepts (e.g., using "spider-man" as an example). Beside of conditional information  $c_{sp}$ , to avoid overfitting during tuning, we use a None condition  $c_N$  (" ") to construct anchor images  $x_{ac}$ . The second step is "Scrambled Grid". Firstly, we grid images which generated by specific conditional information  $c_{sp}$  and randomly reorder them to create scrambled images  $x_{sg} = O(SD_{pr}(\epsilon, c_{sp}))$ , where the  $O$  refers to the Scrambled Grid operation and the size of grid in there is  $16 \times 16$ . The reason for using this operation is that such a destruction technique can destroy low-frequency information in the original image while preserving some high-frequency features. The final step in degeneration-tuning process is "Tuning". Firstly, we construct tuning data  $x \in x_{sg} \cup x_{ac}$  using  $x_{sg}$  and  $x_{ac}$ . Next, we utilize this tuning data and their corresponding text conditions  $c \in c_{sp} \cup c_N$  to fine-tune the parameters in  $SD_{pr}$ , which is

represented by  $DT_\theta$ . The training loss  $L_{DT}$  becomes:

$$\begin{aligned} L_{DT} &:= \mathbb{E}_{x,c,t,\epsilon_t \sim N(0,1)} [\|x - DT_\theta(\alpha_t x + \sigma_t \epsilon_t, c, t)\|_2^2], \\ &= \mathbb{E}_{x,c,t,\epsilon_t \sim N(0,1)} [\|\epsilon_t - \epsilon_\theta(x_t, c, t)\|_2^2] \end{aligned} \quad (11)$$

Since the Scrambled Grid operation narrows the distribution between  $x_T$  and  $\hat{x}_T$ , training  $DT_\theta$  to fit this specific data domain becomes easier comparatively. For a single concept in DT, it is sufficient to construct just 800 to 1000  $x_{sg}$  samples and an equal number of  $x_{ac}$  samples. The process of degeneration-tuning requires only a minimal learning rate of  $1e^{-7}$ , which is significantly lower than the learning rate used in traditional stable diffusion fine-tuning ( $1e^{-4}$ ) as reported by [42], and demands a relatively low number of training epochs, approximately 60. The trained components in the SD is the entire U-net framework [44].

## 4 EXPERIMENTS

In this Section, we show the performance of our Degeneration-Tuning (DT) method in shielding various types of concepts. The pre-trained stable diffusion (SD) we utilized is SD-1.5 [46], which has been open-sourced in Huggingface [20]. All of our training experiments were conducted on a single machine equipped with 8-GPU V100 GPUs. The batch size was set to 16, the learning rate was uniformly adjusted to  $1e^{-7}$ , and the training epochs were consistently set at 60. For each individual concept, we set the number of  $X_{sg}$  and  $X_{ac}$  samples at 900 and 1200, respectively. When masking

multiple concepts simultaneously, it is sufficient to simply stack multiple degraded datasets together.

#### 4.1 Effectiveness in Recontextualization

When we aim to shield a concept, our goal is not only to ensure the model can be activated when a single textual concept is inputted, but also to be effective in all contexts containing that concept. In DT, although we only target specific concepts as the model's shield goals, we find that it can still be effectively applied to various contexts containing these textual concepts. Taking the concept of "superman" as an example. After applying DT to the concept word "superman" (DT on "superman"), the model can shield content not only in single-object prompts such as "a photo of superman" but also in prompts containing multiple objects and meanings, like "children play with superman" and "A cup with a superman logo." The detailed results are presented in Figure 5. By comparing the generated images of the original stable diffusion model and the model after DT on "superman", we demonstrate that our DT method is effective in recontextualization.

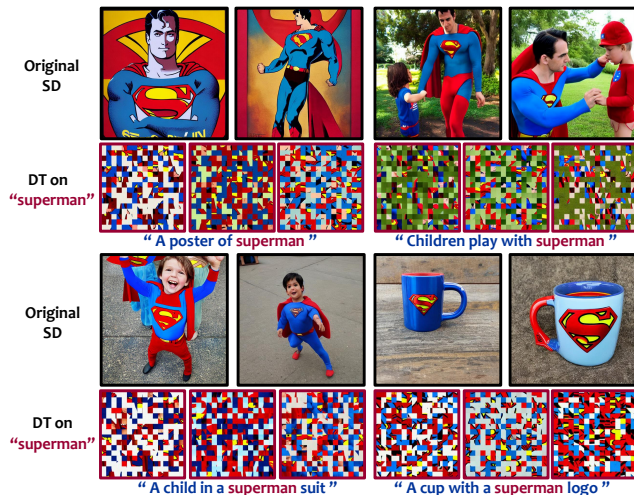


Figure 5: The performance of the SD in different contexts containing the concept "superman" is demonstrated before and after applying DT, using 'superman' as the conditional information.

#### 4.2 Effectiveness in protecting artistic styles

Imitation and plagiarism of artistic styles are the most common copyright problems in generative models. In fact, recent lawsuits against companies such as Stability AI, DeviantArt, and Midjourney highlight the seriousness of this issue [29]. To demonstrate that our degeneration-tuning (DT) method remains effective in masking artistic styles, we applied it to the concepts of "Monet" and "Starry Night." The results, as displayed in Figure 6, indicate that the DT method still excel in shielding the content of artistic styles.

#### 4.3 Effectiveness in various other concepts

In addition to demonstrating the effectiveness of our DT method in shielding specific concepts, in there, taking the concept "spider-man" as an example, we showcase its generative performance on

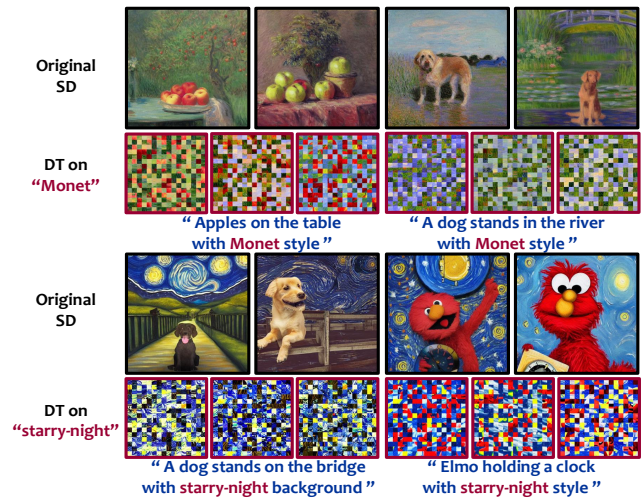


Figure 6: The performance of the SD in generating styles characteristic of 'Monet' and 'starry-night', before and after DT on the textual concepts "Monet" and "starry-night".

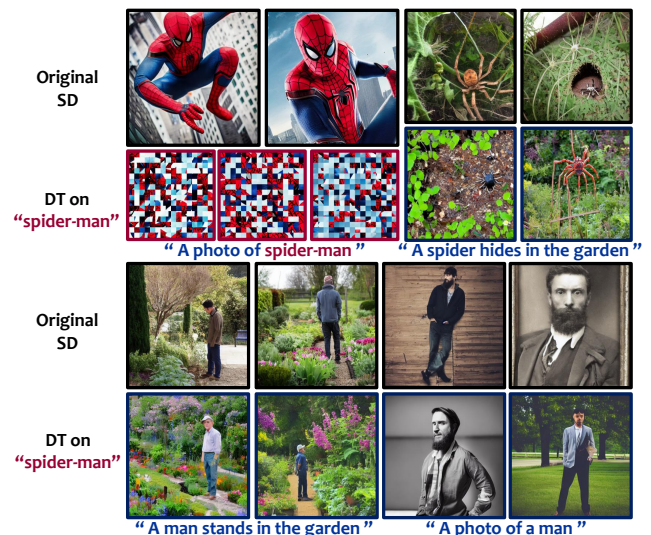


Figure 7: The performance of the SD in concepts related to "spider-man" is shown both before and after DT on the textual concept "spider-man".

non-specific conceptual content. As shown in Figure 7, after DT on textual concept "spider-man" in SD, we can observe that the model does not exhibit significant deviation or degradation in generating content for the concepts of "spider" and "man".

#### 4.4 Effectiveness in Grafting

In addition to the performance of the model after DT, we observe that the  $DT_{\theta}$  model, which has been degeneration-tuned in specific textual concepts, demonstrates a grafting ability. Specifically, when we replace the U-net network of ControlNet with the model  $DT_{\theta}$ , this grafted ControlNet (Con-DT) is able to shield these specific concepts present in the model  $DT_{\theta}$ , even when additional conditional information such as pose and canny information is inputted



along with the textual information. We demonstrate the model’s grafting ability in Figures 1 and 8 for the concepts "spider-man" and "Emma Watson", respectively. In Figure 1, we show that the Con-DT can generate images based on both pose and text information, such as 'robot-man', while effectively shielding the content about 'spider-man', when we graft the U-net after DT on the 'spider-man' concept into a pose-based ControlNet. Similarly, in Figure 8, the performance of the edge-based Con-DT aligns with pose-based Con-DT. The edge-based Con-DT is able to generate content based on canny edge and text information, such as "a photo of Taylor Swift", while still shielding the content about 'Emma Watson'. We believe that the reason for this is that conditional information from different modalities refers to different semantic content, and the information spectra of these semantic content are not the same. Masking the textual concepts does not affect the influence of other modalities’ conditional information on the generative content.

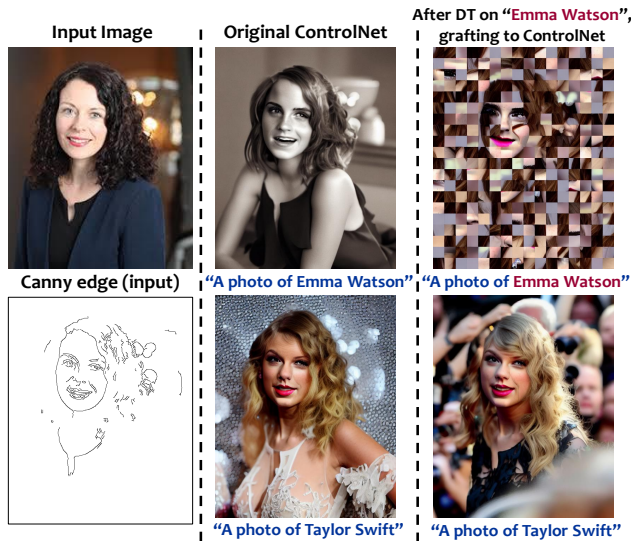


Figure 8: The performance of the edge-based ControlNet in the concept of Emma Watson when given a canny edge input, before and after being grafted from the model that underwent DT on the text "Emma Watson".

### 4.5 Fine-tuning using normal images

In Section 3, we noted that Scrambled Grid (SG) strategy is critical to the effectiveness of the degeneration-tuning method. In there, we show the results of the DT method without SG strategy. Taking the concept "spider-man" as an example, we replaced the data  $x_{sg}$ , which underwent SG, with specific content images, and combined them with images generated from none condition information (" ") to create the tuning dataset  $x$ . The results are shown in Figure 9. Despite elevating the learning rate from  $1e^{-7}$  to  $1e^{-6}$  and extending the training epoch from 100 to 200, the fine-tuned model failed to alter its generative content when supplied with the condition text "spider-man" (first-row). When we increased the learning rate to  $1e^{-5}$ , the model no longer generate images associated with the concept of "spider-man". However it also lost the ability to modulate the generated contents based on textual information (last-two rows). This situation indicates that the model has overfitted to the tuning

dataset  $x$ , emphasizing the difficulties faced by the diffusion model in learning low-frequency signals as compared to high-frequency.

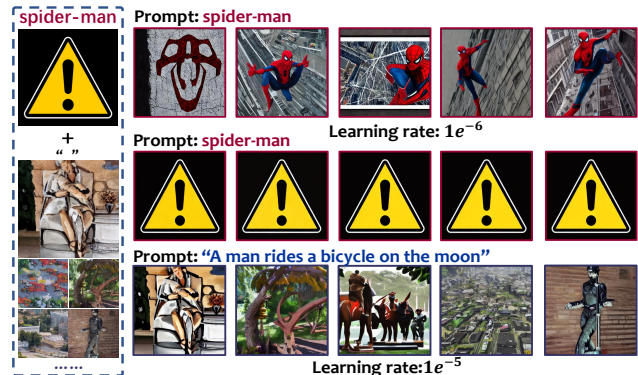


Figure 9: The performance of DT on the text "spider-man" with different learning rates, without using the Scrambled Grid operation.

## 5 EVALUATION

In this section, we adopt the evaluation metrics used in previous generative works [16–18, 43], using FID [10] and IS [55] scores to assess the image generation quality of the model after degeneration-tuning (DT) in various types of concepts and multiple concepts simultaneously. The FID score measures the distance between the generated image distribution and the target distribution, with a smaller score indicating greater similarity. And IS score assesses image fidelity based on a pre-trained Inception-v3 model [55], with a higher score indicating better quality. Beside of this, we discuss the feasibility and challenges of continual degeneration-tuning, which becomes relevant as the number of copyright-related content continues to increase.

### 5.1 Quantitative analysis

We present a quantitative evaluation of the image quality generated by the models after applying DT to different concepts, using FID and IS scores. To demonstrate the effectiveness of our degeneration-tuning method, we separately evaluate the model’s ability to generate both unwanted conceptual images and other general images. For specific conceptual content, we establish the image distribution generated by the original stable diffusion using these textual concepts as the target distribution. Concurrently, we use the image distribution generated by the model after DT on these textual concepts as the generated distribution. For other general content, we align with previous works [22, 37, 43] utilizing prompts from the COCO-30k validation set [28] to generate images and comparing them to the original COCO-30k distribution.

Firstly, we evaluate the changes in the image generation quality of the stable diffusion model before and after degeneration-tuning on a single concept. The detailed results can be seen in the Table 1, where the C.s.c is "Contents about specific concepts".

When comparing the FID and IS scores for specific concepts, we can find that the content about this concepts, which has been degeneration-tuned in the pre-trained stable diffusion model, cannot be generated. Additionally, the results on COCO-30K suggest that our method has little impact on generated contents outside



DT in	C.s.c		COCO 30K	
	FID ↓	IS ↑	FID ↓	IS ↑
<b>Original SD</b>	\	\	<b>12.61</b>	<b>39.20</b>
"spider-man"	385.38	1.77	12.64	38.77
"superman"	371.64	1.63	12.53	38.94
"Monet"	355.20	1.81	12.60	39.12
"starry-night"	360.02	1.70	12.62	38.73
"Emma Watson"	392.21	1.58	12.58	38.60
"Donald Trump"	381.71	1.60	12.70	39.01
"Joint"	391.54	1.73	13.04	38.25

**Table 1: The FID and IS scores of the model before and after DT for unwanted concepts in the concept-specific content and COCO 30K dataset.**

of the specific concepts. Furthermore, we perform DT on all of this concepts jointly, and the FID and IS scores for this joint DT application are presented in the last row of Table 1. By comparing the result of the model after DT with the original SD, we observe that the number of unwanted concepts in DT does not significantly affect the model’s generation quality on other content. Specifically, the FID score for COCO-30K increased by a mere 0.23 points, and the IS score decreased by 0.62 points, reinforcing the effectiveness of our method without compromising the generation quality of other content.

Method	Venue	FID ↓ / IS ↑	CLIP	R.p.l
SLD[48]	CVPR’23	18.76 / 36.64	0.1594	✗
Erase[7]	Arxiv’23	17.27 / 37.21	0.1586	✓
<b>DT</b>	<b>Our</b>	<b>13.04 / 38.25</b>	<b>0.1572</b>	<b>✓</b>
Original	CVPR’22	12.61 / 39.20	0.1561	✗

**Table 2: Comparing the performance of the DT method with existing content protection methods, SLD [48] and Erase [7].**

In addition to evaluating the efficacy of the DT method, we also compared the performance and differences of the DT method with existing content protection methods, such as SLD [48] and Erase [7]. As shown in the Table ??, the resulting average FID and IS scores, along with CLIP scores, clearly show that the DT method maintains a closer resemblance to the original content in terms of general content generation compared to other methods. Among these concurrent works, only DT and Erase methods effectively prevent the risks of uncontrolled generation in the event of model parameter leakage by modifying the model parameters.

## 5.2 Continual degeneration-tuning

If we consider the training process of stable diffusion from a probabilistic perspective, degeneration-tuning the parameters is tantamount to finding their most probable values given some data  $D_j$ , where the  $D_j = D_{sg} \cup D$ . The  $D_{sg}$  is the data after scrambled grid operation, and  $D$  is the data used to pre-train the SD. The probability equation can be rearranged to:

$$\log p(\theta|D_j) = \log p(D_{sg}|\theta) + \log p(\theta|D) - \log p(D_{sg}), \quad (12)$$

This problem is challenging to address in traditional continual tasks where the original data  $D$  is unavailable. However, for generative

models like SD,  $D$  can be re-generated from the model itself. This is the reason why we need to generate  $x_{ac}$  in DT (in Section 3).

In there, we evaluated the performance of continual degeneration-tuning in unwanted concepts. The detailed FID and IS scores has been shown in Table 3. where the C.s.c is "Contents about specific

DT in	C.s.c		COCO 30K	
	FID ↓	IS ↑	FID ↓	IS ↑
<b>Original SD</b>	\	\	<b>12.61</b>	<b>39.20</b>
"spider-man" ↓	385.38	1.77	12.64	38.77
"superman" ↓	380.12	1.70	13.10	38.02
"Monet" ↓	371.62	1.75	13.63	37.14
"starry-night" ↓	392.48	1.68	14.21	36.58
"Emma Watson" ↓	387.61	1.76	14.87	36.21
"Donald Trump" ↓	390.25	1.71	15.32	35.71
"Joint"	391.54	1.73	13.04	38.25

**Table 3: The FID and IS scores of the model after continual DT in different concepts.**

concepts". From the scores, we observe that although continual DT does not significantly affect the model’s performance in shielding specific conceptual contents, its performance on other prompts is negatively impacted. Compared to the joint results, the final FID score of the model on COCO-30K increased from 13.04 to 15.32, while the IS scores decrease from 38.25 to 35.71. By examining the quality of the images generated by the model after each continual tuning phase, we find that continual training amplified the bias between the generated image quality and the original image, leading to a butterfly effect and resulting in a deterioration of the model’s generated image quality.

## 6 CONCLUSION

In this paper, we conducted experiments to demonstrate that the primary factor influencing the semantic contents generated by stable diffusion is the distance between the initial sampled Gaussian noise and the final diffusion distribution within the training data domain. Building on this insight, we proposed a new method **Degeneration Tuning(DT)** to shield unwanted concepts from the level of stable diffusion’s weights. In addition to showing the performance of our DT method in shielding various types of concepts qualitatively, the quantitative comparison of SD before and after DT indicates that DT method did not significantly affect the generative quality of other contents. Finally, we evaluated the feasibility of continual DT and proposed potential reasons that may lead to a decrease in the quality of generated image by the contunual model.

## ACKNOWLEDGEMENTS

This work has been supported in part by the Zhejiang NSF (LR21F020004), the NSFC (No.62272411), Alibaba-Zhejiang University Joint Research Institute of Frontier Technologies, and Ant Group.

## REFERENCES

- [1] Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Schwag, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. 2023. Extracting training data from diffusion models. *arXiv preprint arXiv:2301.13188* (2023).

- [2] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems* 34 (2021), 8780–8794.
- [3] Eliot Higgins. 2022. Making Fake news using Stable Diffusion. <https://twitter.com/EliotHiggins>.
- [4] Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 12873–12883.
- [5] Kevin Frans, Lisa Soros, and Olaf Witkowski. 2022. Clipdraw: Exploring text-to-drawing synthesis through language-image encoders. *Advances in Neural Information Processing Systems* 35 (2022), 5207–5218.
- [6] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618* (2022).
- [7] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. 2023. Erasing concepts from diffusion models. *arXiv:2303.07345* (2023).
- [8] Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. 2013. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211* (2013).
- [9] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 2022. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626* (2022).
- [10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* 30 (2017).
- [11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.
- [12] Jonathan Ho and Tim Salimans. 2022. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598* (2022).
- [13] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685* (2021).
- [14] Rongjie Huang, Feiyang Chen, Yi Ren, Jinglin Liu, Chenye Cui, and Zhou Zhao. 2021. Multi-singer: Fast multi-singer singing voice vocoder with a large-scale corpus. In *Proceedings of the 29th ACM International Conference on Multimedia*. 3945–3954.
- [15] Rongjie Huang, Chenye Cui, Feiyang Chen, Yi Ren, Jinglin Liu, Zhou Zhao, Baoxing Huai, and Zhefeng Wang. 2022. Singgan: Generative adversarial network for high-fidelity singing voice generation. In *Proceedings of the 30th ACM International Conference on Multimedia*. 2525–2535.
- [16] Rongjie Huang, Max WY Lam, Jun Wang, Dan Su, Dong Yu, Yi Ren, and Zhou Zhao. 2022. FastDiff: A Fast Conditional Diffusion Model for High-Quality Speech Synthesis. *arXiv preprint arXiv:2204.09934* (2022).
- [17] Rongjie Huang, Yi Ren, Jinglin Liu, Chenye Cui, and Zhou Zhao. [n. d.]. GenerSpeech: Towards Style Transfer for Generalizable Out-Of-Domain Text-to-Speech. In *Advances in Neural Information Processing Systems*.
- [18] Rongjie Huang, Zhou Zhao, Huadai Liu, Jinglin Liu, Chenye Cui, and Yi Ren. 2022. Prodiff: Progressive fast diffusion model for high-quality text-to-speech. In *Proceedings of the 30th ACM International Conference on Multimedia*. 2595–2605.
- [19] huggingface. 2022. stable diffusion 1-4 and Negative Prompt. <https://huggingface.co/CompVis/stable-diffusion-v-1-4-original/tree/main>.
- [20] huggingface. 2022. web of huggingface. <https://huggingface.co/>.
- [21] Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4401–4410.
- [22] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. 2022. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276* (2022).
- [23] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. 2022. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2426–2435.
- [24] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- [25] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences* 114, 13 (2017), 3521–3526.
- [26] laion.ai. 2020. LAION dataset. <https://laion.ai/>.
- [27] Juncheng Li, Xin He, Longhui Wei, Long Qian, Linchao Zhu, Lingxi Xie, Yueting Zhuang, Qi Tian, and Siliang Tang. 2022. Fine-grained semantically aligned vision-language pre-training. *Advances in neural information processing systems* 35 (2022), 7290–7303.
- [28] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 740–755.
- [29] Stable Diffusion litigation. 2022. negative prompt. <https://stablediffusionlitigation.com/>.
- [30] Zhiheng Liu, Ruili Feng, Kai Zhu, Yifei Zhang, Kecheng Zheng, Yu Liu, Deli Zhao, Jingren Zhou, and Yang Cao. 2023. Cones: Concept neurons in diffusion models for customized generation. *arXiv preprint arXiv:2303.05125* (2023).
- [31] Zhiheng Liu, Yifei Zhang, Yujun Shen, Kecheng Zheng, Kai Zhu, Ruili Feng, Yu Liu, Deli Zhao, Jingren Zhou, and Yang Cao. 2023. Cones 2: Customizable Image Synthesis with Multiple Subjects. *arXiv preprint arXiv:2305.19327* (2023).
- [32] Michael McCloskey and Neal J Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*. Vol. 24. Elsevier, 109–165.
- [33] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaoju Qie. 2023. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453* (2023).
- [34] Zixuan Ni, Haizhou Shi, Siliang Tang, Longhui Wei, Qi Tian, and Yueting Zhuang. 2021. Revisiting Catastrophic Forgetting in Class Incremental Learning. *arXiv:arXiv:2107.12308*
- [35] Zixuan Ni, Longhui Wei, Siliang Tang, Yueting Zhuang, and Qi Tian. 2023. Continual Vision-Language Representation Learning with Off-Diagonal Information. *arXiv preprint arXiv:2305.07437* (2023).
- [36] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2021. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741* (2021).
- [37] openai. 2022. DALLE2. <https://github.com/openai/dalle-2-preview/blob/main/system-card.md>.
- [38] Vinay V Ramasesh, Ethan Dyer, and Maithra Raghu. 2020. Anatomy of catastrophic forgetting: Hidden representations and task semantics. *arXiv preprint arXiv:2007.07400* (2020).
- [39] Javier Rando, Daniel Paleka, David Lindner, Lennard Heim, and Florian Tramèr. 2022. Red-Teaming the Stable Diffusion Safety Filter. *arXiv preprint arXiv:2210.04610* (2022).
- [40] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. 2017. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. 2001–2010.
- [41] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. 2014. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*. PMLR, 1278–1286.
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. *arXiv:2112.10752* [cs.CV]
- [43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10684–10695.
- [44] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. 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*. Springer, 234–241.
- [45] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2022. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arXiv:2208.12242* (2022).
- [46] runwayml. 2022. stable-diffusion-v1-5. <https://huggingface.co/runwayml/stable-diffusion-v1-5/tree/main>.
- [47] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* 35 (2022), 36479–36494.
- [48] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. 2023. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 22522–22531.
- [49] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*. PMLR, 2256–2265.
- [50] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502* (2020).
- [51] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. 2021. Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems* 34 (2021), 1415–1428.
- [52] stabilityai. 2022. negative prompt. <https://huggingface.co/spaces/stabilityai/stable-diffusion/discussions/7857>.
- [53] stability.ai. 2022. stable-diffusion-v2-release. <https://stability.ai/blog/stable-diffusion-v2-release>.

- [54] Shikun Sun, Longhui Wei, Junliang Xing, Jia Jia, and Qi Tian. 2023. SDDM: Score-Decomposed Diffusion Models on Manifolds for Unpaired Image-to-Image Translation. (2023).
- [55] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2818–2826.
- [56] Longhui Wei, Lingxi Xie, Wengang Zhou, Houqiang Li, and Qi Tian. 2022. Mvp: Multimodality-guided visual pre-training. In *European Conference on Computer Vision*. Springer, 337–353.
- [57] wv2nw0. 2022. how to remove the safety filter in 5 seconds. [https://www.reddit.com/r/StableDiffusion/comments/wv2nw0/tutorial\\_how\\_to\\_remove\\_the\\_safety\\_filter\\_in\\_5/](https://www.reddit.com/r/StableDiffusion/comments/wv2nw0/tutorial_how_to_remove_the_safety_filter_in_5/).
- [58] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. 2021. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157* (2021).
- [59] Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. 2021. Vector-quantized image modeling with improved VQGAN. *arXiv preprint arXiv:2110.04627* (2021).
- [60] Lvmin Zhang and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543* (2023).



## A ANALYSIS AND ABLATION

In this Appendix, we provide further analysis of the mechanism behind the Degeneration-Tuning (DT) method and experimental results. Furthermore, we showcase various ablation experiments based on the DT proposed in the Section Methods.

### A.1 The impact of the degeneration-tuning on the original model parameters

In Section Methods, we mentioned that the DT method can fit the data processed by Scrambled Grid with a very small learning rate (1e-7). However, how much impact does it have on the original model parameters? Here we present an experiment to demonstrate it. Taking the concept "spider-man" as an example. Assuming the original pre-trained stable diffusion parameters are  $\theta_{ori}$ , and the model parameters after DT are  $\theta_{DT}$ . By performing a linear transformation, we obtain fused models  $\theta_f$  under different hyper-parameter  $\lambda$ , where the  $\theta_f = \lambda\theta_{ori} + (1 - \lambda)\theta_{DT}$ . Figure 10 shows the masking effect of the model  $\theta_f$  on the prompt "spider-man stand in the garden" under different hyper-parameter  $\lambda$ . From the experiments,

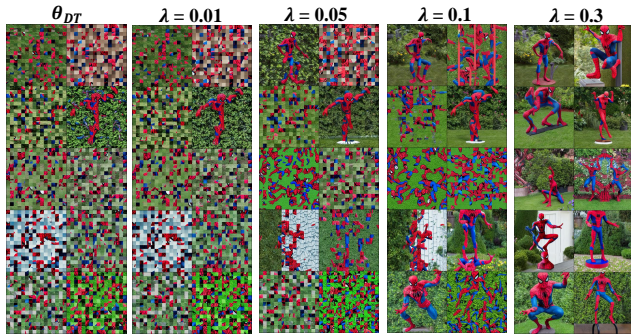


Figure 10: The performance of fusion parameters generated by different hyper-parameter  $\lambda$  on the prompt "spider-man stands in the garden".

we can find that the shielding effect of the model on the textual concept "spider-man" starts to weaken when the hyper-parameter  $\lambda$  changes from 0.01 to 0.05. When the hyper-parameter reaches 0.1, the model becomes sluggish towards the concept of spider-man. All of this demonstrates that although the influence of the DT on the original model parameters is small, its effect on the semantic contents of the generative image is significant.

### A.2 Applying DT to different modules of the model

In Section Methods, we state that the DT method was used to adjust the parameters within the U-Net framework of stable diffusion. However, the U-Net framework includes the cross-attention modules and the resblock modules. Here, taking the concept "spider-man" as an example, we discuss the performance when applying DT only to the cross-attention modules or the resblock modules.

**A.2.1 Just applying DT to the cross-attention modules.** After adjusting only the parameters of the cross-attention modules within

the U-Net framework using DT with the textual concept of "spider-man", we show the performance of this model in the prompt "A man stands in the garden" without the concept "spider-man" in Figure 11. We can clearly see that the content of these images contains textual information, but its saturation and contrast are very high, which makes the entire output appear unrealistic.



Figure 11: The performance of the model, which applies DT only to the cross-attention modules, on the prompt 'A man stands in the garden'.

**A.2.2 Just applying DT to resblock modules.** After adjusting only the parameters of the resblock modules within the U-Net framework using DT with the textual concept of "spider-man", we show the performance of this model in the prompt "A man stands in the garden" without the concept "spider-man" in Figure 12. We can clearly see that the content of these images is always different from textual information.



Figure 12: The performance of the model, which applies DT only to the resblock modules, on the prompt 'A man stands in the garden'.

### A.3 Direct fine-tuning using pure color images

In Section Experiments, we show the performance of DT without the Scrambled Grid operation. In this section, we supplement the results of using images with a single low-frequency signals (black, white and gray) instead of scrambled images in Figure 14. From the results, we can see that regardless of which solid color image is used as a substitute for the text concept "spider-man", the model still generates images related to the concept "spider-man" when it receives conditional information about "spider-man". All of this once again indicates that the model's learning of low-frequency information is more difficult compared to high-frequency information.



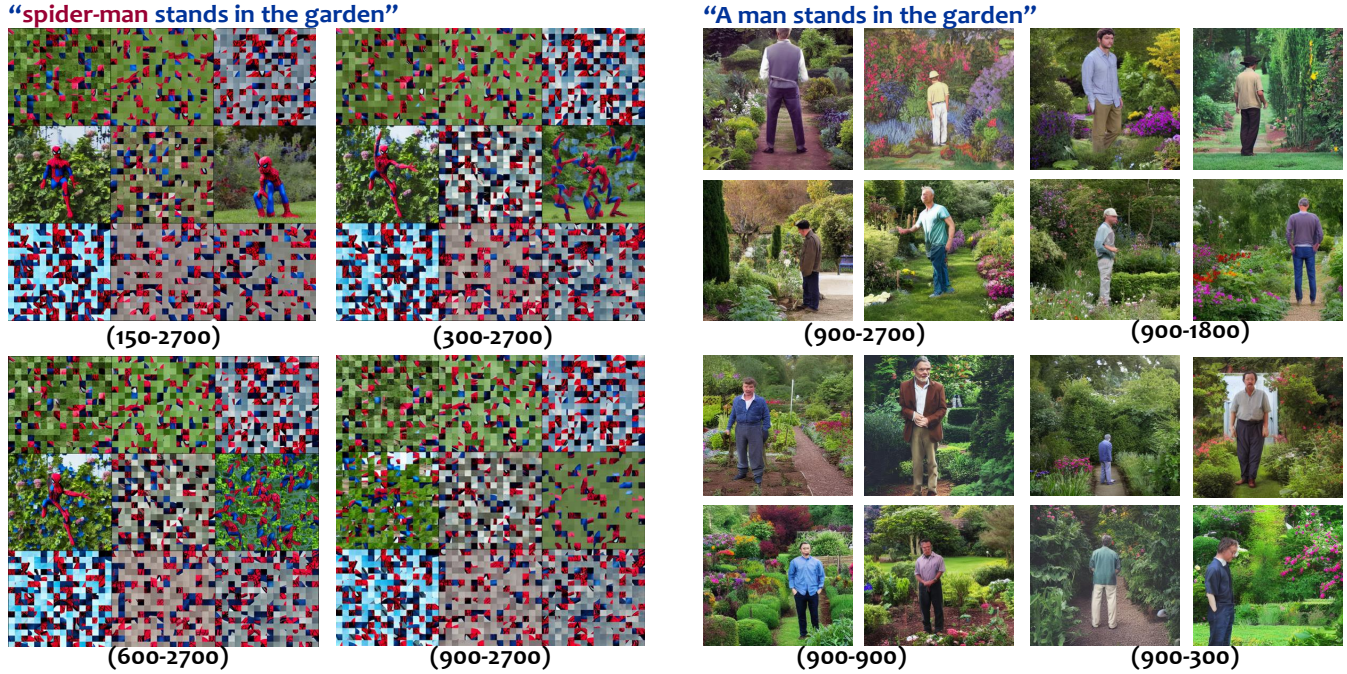


Figure 13: Taking the concept "spider-man" as an example, we show the impact of different ratios and numbers of  $x_{sg}$  and  $x_{ac}$  images on DT performance.

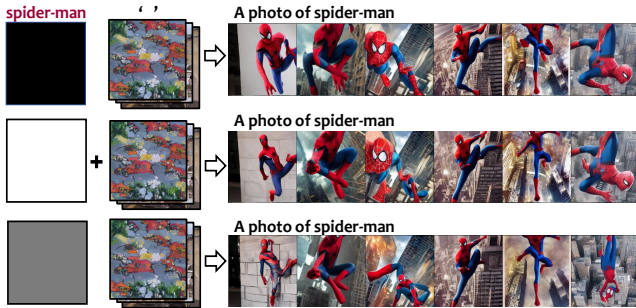


Figure 14: The result of the model after fine-tuning on different replacement dataset.

#### A.4 The proportion of anchor images in Degeneration-Tuning

In Section Experiments, we set the proportion of  $x_{sg}$  to  $x_{ac}$  to be approximately 1:1. In this section, taking the concept "spider-man" as an example, we show the impact of different ratios and numbers of  $x_{sg}$  and  $x_{ac}$  images on DT performance. From Figure 13, we can find that when the number of  $x_{sg}$  increases from 150 to 900, its shielding effect on the concept "spider-man" varies. As the number of  $x_{sg}$  increases, its effectiveness improves progressively. When we fix the number of  $x_{sg}$  at 900 and adjust the quantity of  $x_{ac}$ , as shown in Figure 13, we find that for individual concepts, the performance of DT on other contents does not show significant changes. However, during the actual experimental process, we still set the number of  $x_{ac}$  at 900~1200 to reduce the impact of DT on other contents.

#### A.5 The size of Scrambled Grid

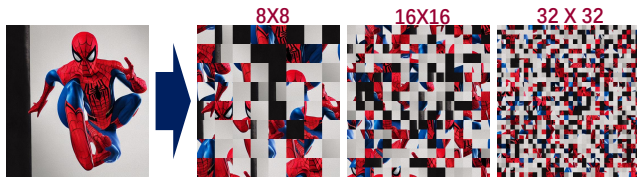
In Section Methods, we stated that the size of Scrambled Grid in experiments were set to  $16 \times 16$ . In this Section, we show the quantitative results of the model after DT with different size of Scrambled Grid on concept "spider-man". Figure 15 shows the demonstration of Scrambled Grid with different grid size. And Table 4 shows the FID and IS scores of the model after DT with different size of Scrambled Grid on concept "spider-man". Based on the results, the influence of DT methods with different grid scales on the model's generation quality is not significant. However, the  $8 \times 8$  size visually retains too much image content, while the  $32 \times 32$  size destroys too much high-frequency information from the original content. Therefore, in practical experiments, we chose  $16 \times 16$ .

The size of Scrambled Grid	C.s.c		COCO 30K	
	FID	IS	FID	IS
<b>Original SD</b>	\	\	<b>12.61</b>	<b>39.20</b>
8X8	357.21	1.80	12.63	38.69
16X16	<b>385.38</b>	<b>1.77</b>	<b>12.64</b>	<b>38.77</b>
32X32	384.62	1.72	12.59	38.72

Table 4: The FID and IS scores of the model after DT with different size of Scrambled Grid on concept "spider-man".

#### A.6 Another Performance of DT in Grafting

We demonstrate the another performance model's grafting ability in Figures 16 for the concepts "Donald Trump". When we graft the



**Figure 15: Demonstration of Scrambled Grid with different grid size.**

model after DT on concept "Donald Trump" into a pose-based ControlNet as shown in Figure 16, we show that Con-DT can generate images based on both pose and text information, such as "a photo of Obama", while still effectively shielding the content about "Donald Trump".

### **A.7 The performance of the model on COCO 30K prompts after DT in multiple concepts.**

In Section Evaluation, we compared the quantitative results of the model before and after DT on COCO 30K prompts. In this Section,

we show the generative images by the model before and after DT on COCO 30K prompts qualitatively. Figure 17 shows the performance of the original stable diffusion on COCO 30K prompts. And Figure 18 shows the performance of the stable diffusion on COCO 30K prompts after DT in multiple concepts.

### **A.8 The performance of the model on COCO 30K prompts after continual DT on multi concepts.**

In Section 5.2, we compared the quantitative results of the model before and after continual DT on COCO 30K prompts. In this Section, we show the generative images by the model on some COCO 30K prompts after continual DT on multiple concepts. Figure 19 and 20 show the performance of the this model on some COCO 30K prompts.





Figure 16: The performance of the pose-based ControlNet in the content of Donald Trump when given a pose image input, before and after being grafted from the model that underwent DT on the textual concept "Donald Trump".



Figure 17: The images generated by original stable diffusion on some COCO 30K prompts.



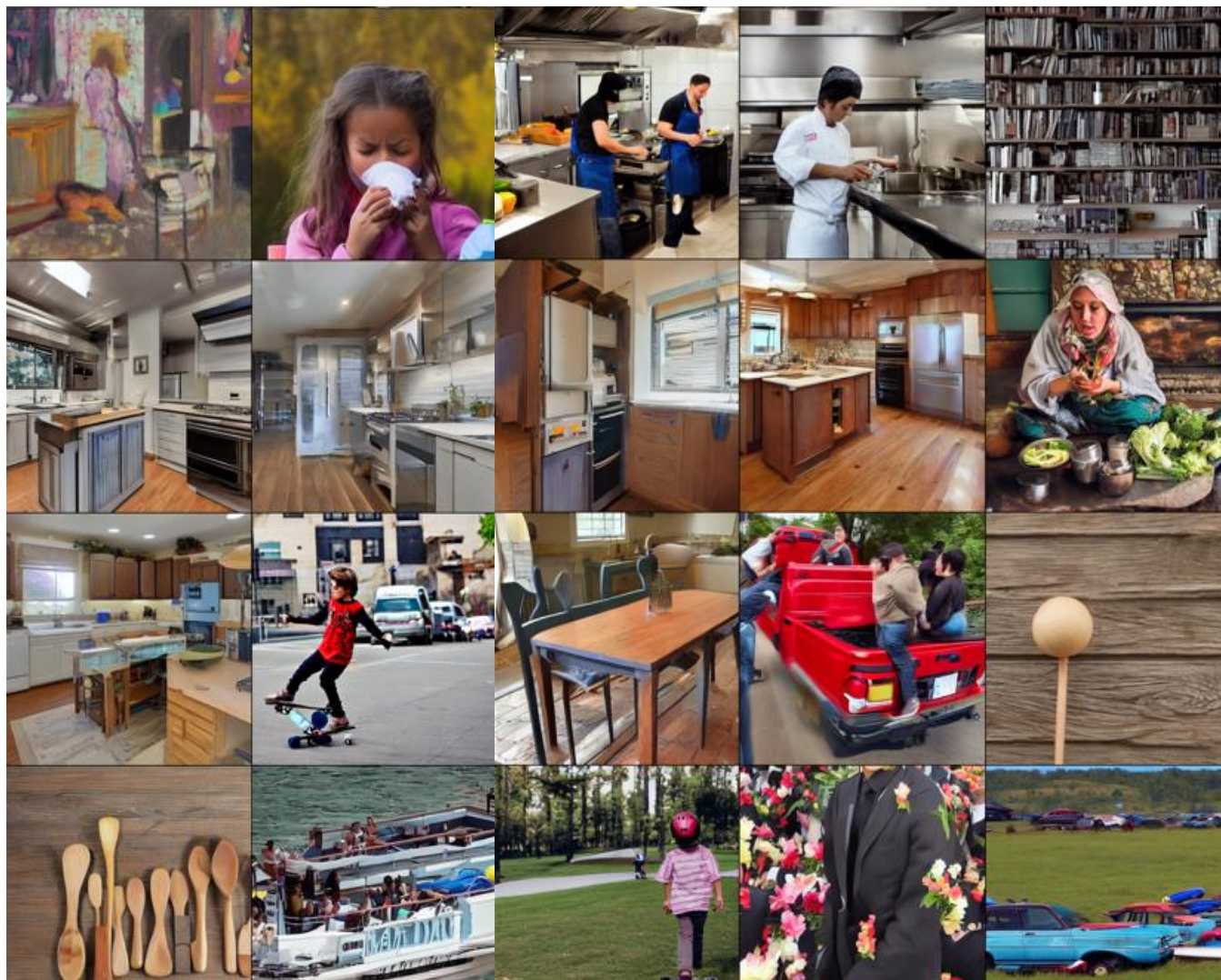


Figure 18: The images generated by the model on some COCO 30K prompts after DT on multiple concepts.

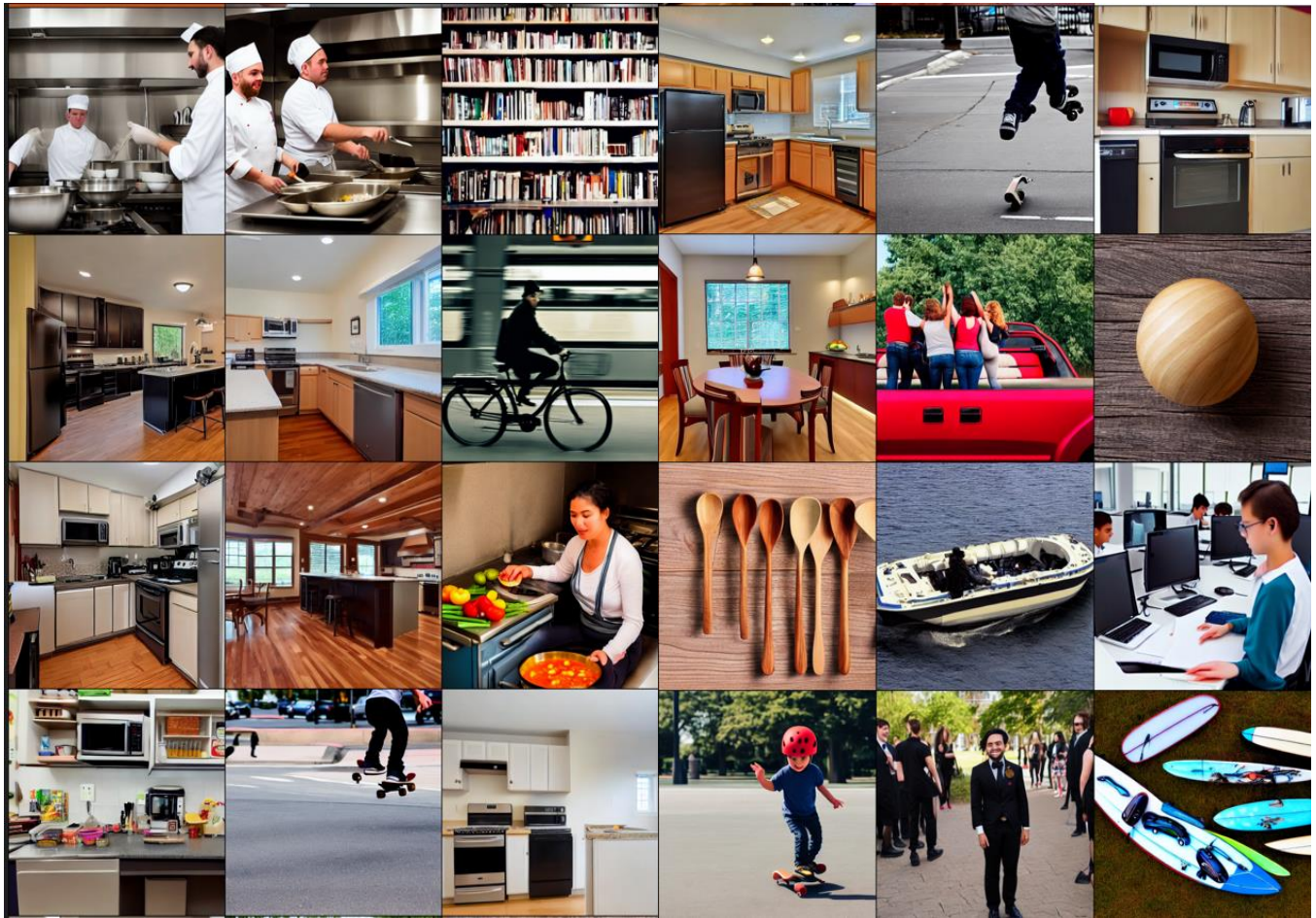


Figure 19: The generative images by the model on some COCO 30K prompts after continual DT on multiple concepts.



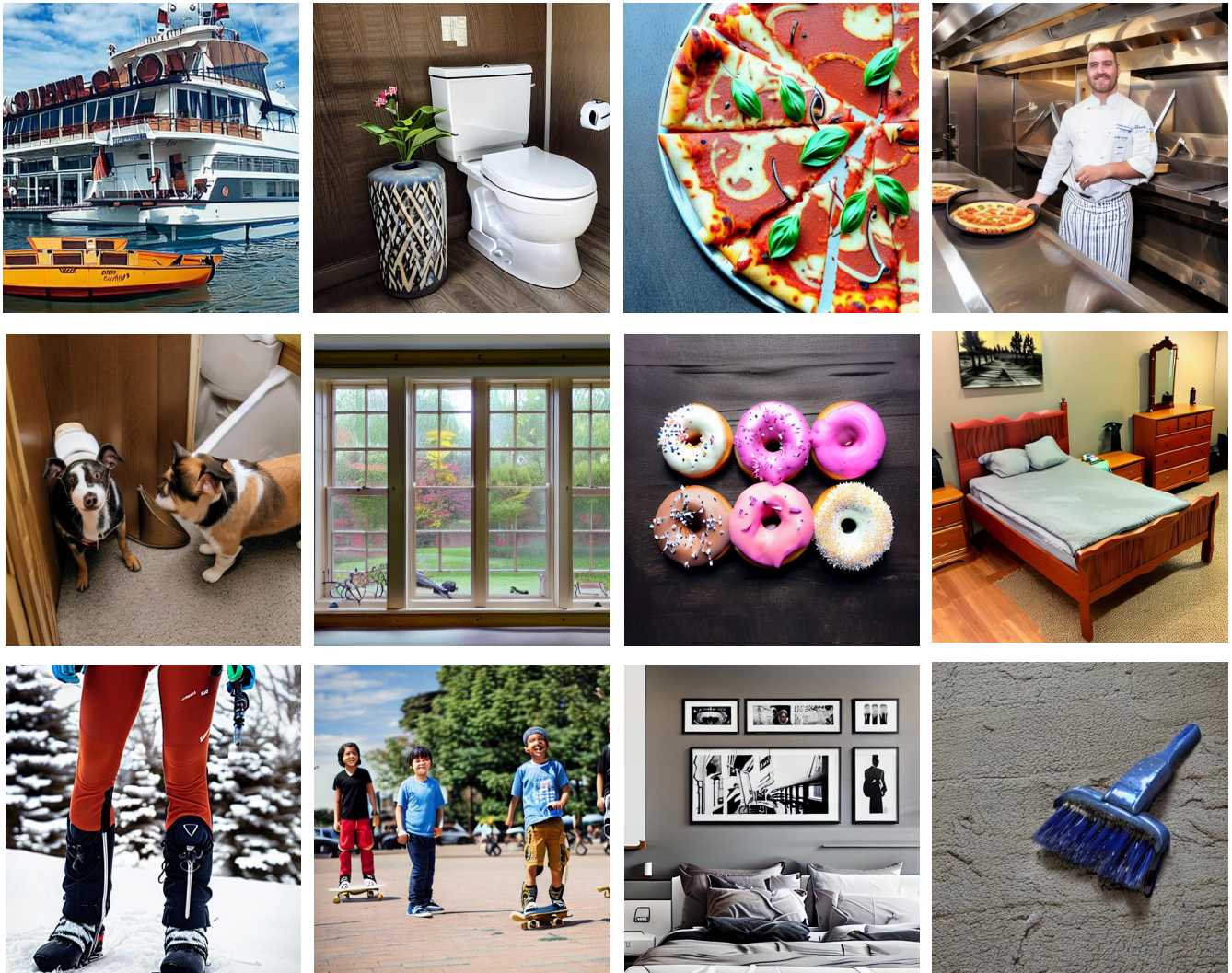


Figure 20: The generative images by the model on some COCO 30K prompts after continual DT on multiple concepts.