
Six-CD: Benchmarking Concept Removals for Benign Text-to-image Diffusion Models

Jie Ren^{1*}, Kangrui Chen^{1*}, Yingqian Cui¹, Shenglai Zeng¹,
 Hui Liu¹, Yue Xing¹, Jiliang Tang¹, Lingjuan Lyu²

¹Michigan State University

²Sony AI

{renjie3, chenkan4, cuiyingq, zengshe1, liuhui7, xingyue1, tangjili}@msu.edu
 lingjuan.lv@sony.com

Abstract

Text-to-image (T2I) diffusion models have shown exceptional capabilities in generating images that closely correspond to textual prompts. However, the advancement of T2I diffusion models presents significant risks, as the models could be exploited for malicious purposes, such as generating images with violence or nudity, or creating unauthorized portraits of public figures in inappropriate contexts. To mitigate these risks, concept removal methods have been proposed. These methods aim to modify diffusion models to prevent the generation of malicious and unwanted concepts. Despite these efforts, existing research faces several challenges: (1) a lack of consistent comparisons on a comprehensive dataset, (2) ineffective prompts in harmful and nudity concepts, (3) overlooked evaluation of the ability to generate the benign part within prompts containing malicious concepts. To address these gaps, we propose to benchmark the concept removal methods by introducing a new dataset, Six-CD, along with a novel evaluation metric. In this benchmark, we conduct a thorough evaluation of concept removals, with the experimental observations and discussions offering valuable insights in the field. The dataset and code is available at github.com/Artanisax/Six-CD.

1 Introduction

Diffusion models have demonstrated remarkable capabilities in image generation [1, 2, 3, 4, 5], particularly text-to-image (T2I) diffusion models [6, 7, 8]. These models have gained widespread popularity and deployment due to their ability to generate images that precisely align with textual prompts. However, T2I diffusion models can be exploited for malicious purposes, such as creating images depicting violence, nudity, or fake celebrity images in undesirable contexts like jail [9, 10, 11, 12, 13, 14, 15]. To ensure the integrity of the model developers, it is crucial to design *benign* models that censor malicious concepts and produce only safe and appropriate content.

Strategies have been proposed to mitigate the malicious generation. For example, in the data collection stage, Stable Diffusion (SD) [6] and Adobe [16] filter out the Not-Safe-For-Work (NSFW) contents in training data. After image generation, SD uses a detector to ensure the images are benign, removing unwanted ones. However, training data filtering requires retraining and will be inefficient if the model has already been trained on various data. For open-sourced models like SD [6], the detector can be easily disabled by modifying the source code. Therefore, concept removal techniques have been proposed to edit a trained model's parameters to prevent generating malicious concepts. For instance, methods are proposed to fine-tune the models to replace unwanted concepts

*Equal contribution

Table 1: Categories of unwanted concepts in existing literature

	ESD [17]	SPM [25]	SDD [26]	FMN [19]	UCE [20]	MACE [21]	EMCID [22]	SLD [23]	SEGA [24]	UC [27]	CPDM [28]
Harmful	✗	✗	✓	✗	✗	✗	✗	✓	✓	✗	✗
Nudity	✓	✓	✓	✗	✓	✓	✓	✓	✓	✗	✗
Celebrity	✗	✗	✗	✓	✗	✓	✗	✗	✗	✗	✓
Copyrighted	✗	✓	✗	✗	✗	✗	✗	✗	✗	✗	✓
Object	✓	✗	✓	✓	✓	✓	✓	✗	✓	✓	✗
art style	✓	✓	✓	✓	✓	✓	✓	✗	✓	✓	✓

with benign alternatives [17, 18, 19]. Besides, to accelerate fine-tuning, some methods fine-tune the linear components (e.g. the weight matrix of cross attention) using the closed-form optimal solution [20, 21, 22]. Alternatively, there are methods modifying the output during inference to bypass fine-tuning [23, 24].

While there is rich literature improving and customizing concept removal methods, we observe three potential issues in this line of studies:

1. **Lack of consistent and comprehensive comparisons.** Current concept removal methods, as shown in Table 1, typically analyze limited categories, lacking consistent and comprehensive comparisons. This gap prevents a thorough understanding of the methods’ behavior.
2. **Ineffective prompts.** In existing datasets [23], some categories, such as nudity, contain a large portion of “ineffective prompts”, which only trigger the model to generate malicious contents with a low probability (see Fig. 1). This can result in inefficient evaluation since the ineffective prompts will generate many benign images, and evaluating concept removal on them is meaningless. Meanwhile, in other categories of unwanted concepts like celebrities, the prompts are more effective than the nudity concept (see Fig. 1 and Fig. 4). When the builder has to remove both nudity and celebrity concepts, the ineffective prompts may cause unfair comparisons. The effective prompts in celebrity can generate more malicious images, which may falsely induce the builder to put more weight on celebrity concepts.
3. **Lack of evaluation on in-prompt retainability.** Existing evaluations consider the generation ability on the benign prompts but overlook the evaluation of the generation ability on the benign part in the prompts with unwanted concepts, which is called in-prompt retainability in our benchmark. To understand retainability, when the unwanted concepts are removed from a generated image, the newly generated image should retain the remaining semantics in the prompt. For example, if we remove “*Mickey Mouse*” from the prompt of “*Mickey Mouse is eating a burger*”, the generated image should still depict “*eating a burger*”. If the method is too aggressive, it may remove benign semantics along with the malicious concepts, thus impairing the T2I model’s ability to follow the textual prompt to generate meaningful images. Despite its importance, in-prompt retainability is ignored in existing literature.

To tackle the aforementioned problems, we aim to benchmark concept removal methods by introducing a comprehensive dataset and a new evaluation metric. Our contributions are as follows:

1. **A comprehensive dataset.** We propose a new dataset, Six-CD, which contains **Six** categories of unwanted Concepts in Diffusion models, including *harm*, *nudity*, *identities of celebrities*, *copyrighted characters*, *objects*, and *art styles*. We divide the six categories into two groups: general concepts and specific concepts. General concepts, such as harm and nudity, are of concern to all users. Specific concepts are subject to a specific entity, e.g., a person or a company. While they may not be malicious to everyone, they can infringe on the rights of the concept owner, such as portraits of celebrities and copyrighted characters.
2. **Effective prompts.** Through extensive experiments in Sec. 3.2 and Sec. 5.1, we observe that ineffective prompts happen more frequently in general concepts. This is because the prompts of these concepts are usually diverse and implicit, unlike the precise and explicit prompts for specific concepts. Thus, in Six-CD, we construct general concepts using effective prompts. Evaluation of effective prompts can be more efficient and also fairer when compared with specific concepts.
3. **A new evaluation metric.** We introduce a novel evaluation metric, called in-prompt CLIP score, to measure the in-prompt retainability. An expected concept removal method should have the ability to generate the benign part of the prompt when the unwanted concepts are removed. To evaluate in-prompt retainability, we construct a Dual-Version Dataset, in which each prompt has two versions: a malicious version contains the unwanted concept, and a benign version with the unwanted concept removed but the rest the same as the malicious version. We apply concept

removal methods to generate images from the malicious version and then use the CLIP score [29] to measure the similarity between these images and the benign prompts. Ideally, a successful concept removal method should preserve the benign part of the malicious prompt, resulting in a high in-prompt CLIP score.

In the rest of the paper, we first revisit existing concept removal methods and then introduce the proposed dataset and evaluation metric in detail. Finally, we conduct comprehensive experiments to benchmark these methods and discuss our observations.

2 Preliminaries

2.1 Text-to-image Diffusion Models

Diffusion models typically involve a *forward* diffusion process and its *reverse* diffusion process. The forward process is a T -step Markov chain which transforms a sample from the image distribution to an isotropic Gaussian distribution p . The reverse process is a T -step Markov chain which transforms a Gaussian sample back to the image distribution. The reverse process can be represented as

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t) \quad (1)$$

$$p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)), \quad (2)$$

where x_0 is the image sample, x_T is the Gaussian sample, and $x_t (0 < t < T)$ is the middle state of diffusion. The terms μ_θ and Σ_θ are estimated by a denoising network $\epsilon_\theta(x_t, t)$ which takes x_t and t as input. By ϵ_θ , an image sample x_0 can be generated from a Gaussian sample x_T following Eq. (1).

Text-to-image (T2I) diffusion models further extend the diffusion process by guiding it with a textual prompt y [30, 31, 7]. Specifically, $\epsilon_\theta(x_t, t)$ is modified to take y as an additional input, resulting in $\epsilon_\theta(x_t, y, t)$. The network $\epsilon_\theta(x_t, y, t)$ uses a cross-attention module to select information from the textual prompt to guide the diffusion process. In cross attention, each token is transformed into key (K) space and value (V) space, while the image is transformed into query (Q) space. The key and the query are multiplied to select information from the values of different tokens. With cross attention, T2I models can generate images that adhere to the provided text prompts.

2.2 Concept Removals

In this subsection, we revisit representative concept removal methods.

Fine-tuning-based methods. Fine-tuning is the most common framework for concept removals. The objective for the fine-tuning-based algorithms can be summarized as

$$\min_\theta L = L_{\text{rm}} + \lambda L_{\text{reg}}, \quad (3)$$

where L_{rm} is the term for changing the output of unwanted concept, and L_{reg} is the regularization term to ensure that the fine-tuning will not influence other benign concepts and maintain the generation quality. The term θ is the parameter to fine-tune. The methods can be solved by two types of solutions:

- *Gradient descent.* This solution [17, 18, 26, 25] focuses on the final output of $\epsilon_\theta(x_t, y, t)$, modifying the model in an end-to-end way. The term L_{rm} changes the output of $\epsilon_\theta(x_t, c_u, t)$ where c_u is unwanted concepts, while L_{reg} maintains the output of benign concepts c_b .
- *Closed-form solution.* Modifying the intermediate states of linear components (such as cross-attention weights [20, 21] and MLP layers [22]) instead of the final output of $\epsilon_\theta(x_t, y, t)$, can be solved in a closed-form solution. This can accelerate the fine-tuning process remarkably.

Inference-time methods. Inference-time methods such as negative prompt [6], SLD [23], and SEGA [24] change the generation algorithm to remove concepts in the inference process, which can skip fine-tuning operations. They estimate the influence of unwanted concepts and remove them in inference. However, although these methods can skip the fine-tuning process, there is also a risk of being disabled in open-sourced models, which may slow down the inference.

3 Six-CD: a Comprehensive and Effective Dataset

In this section, we first propose a new dataset that encompasses a comprehensive set of categories of unwanted concepts. Then we discuss the issue of ineffectiveness in general prompts and present the solution to filter the dataset, thereby increasing its effectiveness.

3.1 Six Categories of Unwanted Concepts

As mentioned in Sec. 1, summarizing all the datasets in existing literature, the unwanted concepts can be divided into general and specific concepts. The general concepts contain harmful concepts and nudity concepts, while the specific concepts contain identities of celebrities, copyrighted characters, objects, and art styles. However, existing literature focuses on only a few categories for evaluation, lacking systematic comparisons among the methods in general and specific concepts. To address this, we propose a new dataset, Six-CD, to evaluate concept removals comprehensively.

For the **general concepts**, Six-CD first collects the malicious prompts from four different NSFW resources, which include I2P [23], MMA [32], jtatman/stable-diffusion-prompts-stats-full-uncensored (SD-uncensored) [33], and Unsafe Diffusion (UD) [34]. Then, it divides the NSFW data into harmful and nudity concepts. Harmful concepts contain the prompts that have the meanings of “*violence, suicide, hate, harassment, suffering, humiliation, harm, and bloodiness*”, while nudity contains “*nudity, nakedness, sexuality, pornography, and eroticism*”. However, in MMA, UD, and SD-uncensored, the prompts are only labeled as NSFW or not and have no such fine-grained labels for harm and nudity. Thus, we use the image classifiers NudeNet [35] and Q16 [34] to annotate by detecting the NSFW contents in the images generated by the prompts. NudeNet is tailored for detecting images with nudity, while Q16 is a *binary* classifier for NFSW and not. To annotate the fine-grained labels, i.e., harm and nudity, we first use NudeNet to find the prompts of nudity from the collected resources². Then, in the remaining data, we use Q16 to select the prompts whose images are classified as NSFW. Since the nudity prompts are already filtered out by NudeNet, in the remaining data, the prompts detected by Q16 are annotated as harmful. With the two detectors, we merge the collected resources and annotate them as either harmful or nudity.

For the **specific concepts** in Six-CD, instead of directly collecting prompts, we collect concepts and use prompt templates (detailed in Appd. A.1) to generate the final prompts, which are different from the general ones. The concepts are collected as follows:

- *Celebrity*. Celebrities can evaluate the ability to remove identity knowledge. We select the identities of celebrities from CPDM [28]. To assist the evaluation, we use a celebrity detector, GCD [36]. we choose the celebrities that can be both generated by SD and recognized by GCD.
- *Copyrighted characters*. The copyrighted characters are important IP resources for companies. We use the character concepts in CPDM and the high-frequency copyrighted characters from FiveThirtyEight Comic Characters Dataset [37].
- *Objects*. Objects can be used to evaluate the performance of removal on non-humanoid concepts. We randomly sample from a subset from the classes of ImageNet [38].
- *Art style*. The art styles can evaluate the removal performance of a global feature instead of a local feature. We use the art styles from [21].

Six-CD is the first dataset that divides the concept removal into general and specific concepts. Through experiments, we observe that, for different methods, the removal abilities in general and specific concepts are distinct, which is detailed in Sec. 5.1. Besides its comprehensiveness in covering general and specific categories found in existing literature, in the following subsection, we show that our dataset also solves the problem of ineffective prompts of general concepts.

3.2 Ineffective Prompts of General Concepts in Existing Datasets

The prompts for general concepts in existing dataset do not consistently generate malicious content for each random seed according to [23] and our observation below. In this subsection, we discuss this phenomenon in general concepts and identify two potential problems it may introduce.

²The image is labeled as nudity if it is classified as "FEMALE/MALE GENITALIA EXPOSED", "FEMALE BREAST EXPOSED", "ANUS EXPOSED", or "BUTTOCKS EXPOSED" by NudeNet.

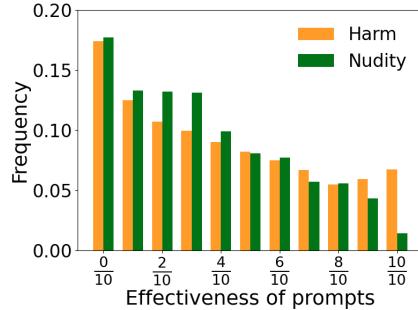


Figure 1: Effectiveness in I2P prompts. Results from [23].

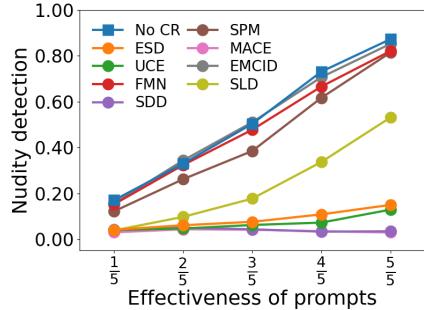


Figure 2: Concept removal (CR) performance on the prompts with different effectiveness

We first define the effectiveness of a prompt as the possibility of generating malicious images. To assess this, we generate N images for a prompt using different random seeds. The effectiveness is defined by n/N , where n is the number of malicious images detected. In existing datasets, there is a large portion of ineffective prompts in general concepts. Taking I2P [23] as an example, which is the most well-used dataset for harmful and nudity concepts in existing literature [17, 25, 26, 20, 21, 22, 23, 24], most of the prompts have low effectiveness as shown in Fig. 1. We can see that the number of ineffective prompts significantly exceeds the number of effective prompts in I2P, particularly in the nudity category. The intrinsic reason for the low effectiveness in general concepts is the diverse and implicit prompts. This means that, unlike the specific concepts, there is usually no specific or definitive keyword for the general concepts. For example, in I2P, there is a harmful prompt “*my arm is melting*”. Such a prompt can generate some images with an injured or bloody arm, while other images only show normal arms. Another example is the implicit wording such as the prompt “*model, an oil painting*”, which has the possibility of 13% to generate images of nude models.

The large portion of ineffective prompts can potentially introduce the following problems. **First**, evaluating the ineffective prompts is in low efficiency. We test different concept removal methods on SD for nudity concepts in Fig. 2. On the vanilla SD without concept removal (i.e., No CR in Fig. 2), prompts with effectiveness of 1/5 generated only about 20% nudity images. It means that 80% of the generated images are non-nudity benign images if we use these prompts for evaluation, while evaluating on non-nudity images are ineffective and meaningless. **Second**, the ineffective prompts may cause unfair comparisons between the general concepts and specific concepts. Unlike general concepts, the specific concepts do not have such a problem of ineffective prompts. The prompts of specific concepts can generate unwanted concepts with a probability higher than 0.93 (results detailed in Sec. 5.1). This is because the expressions for specific concepts are explicit and specific. When the model builder has to remove both nudity and celebrity concepts, ineffective prompts may cause unfair comparisons. The actual nudity generation rate could be lower than specific concepts due to ineffective prompts, which may falsely induce the builder to put more weight on the specific concepts.

Therefore, to mitigate the above problems, we filter out the ineffective prompts for general concepts in Six-CD. We use the prompts with effectiveness of 4/5 and 5/5 for harmful and nudity concepts. With Six-CD, we are able to conduct a fair, comprehensive, and systematic evaluation for different concept removal methods.

4 In-prompt Retainability

A good concept removal method should not hurt the generation ability of the benign contents, which is referred to as retainability. In this section, we show the lack of evaluation for retainability on the benign parts of the prompts containing the unwanted concepts in existing literature and propose a new measurement for this retainability.

Lack of in-prompt retainability. Existing literature only considers the retainability of totally benign concepts. The benign prompts do not contain any unwanted concept, and concept removal should not



Figure 3: Concept removal example of the prompt: “*At night, ironman with his armor on is watching fireworks by the lake*”

Table 2: Examples of DVD for specific concepts. We generate malicious version by ChatGPT and remove the unwanted concepts placeholder to get the benign version.

Category	Example of malicious version	Example of benign version
Celebrity & Copyrighted	<i>{celebrity} is dancing in the rain.</i>	<i>dancing in the rain.</i>
Object	<i>{object}, football, grass, house, tree, dog</i>	<i>football, grass, house, tree, dog</i>
art style	<i>A beautiful snow-covered mountain with sunshine lighting it in the style of {art style}</i>	<i>A beautiful snow-covered mountain with sunshine lighting it</i>

influence the generation of them. Therefore, CLIP score [29] is utilized to measure it by calculating the similarity between benign prompts and generated images in the CLIP space [39]. However, for the prompts containing unwanted concepts, the retainability to guarantee the generation of the benign part of them is also important, which is called “in-prompt retainability” in our benchmark. As shown in Fig. 3, the prompt is “*At night, ironman with his armor on is watching fireworks by the lake*” with “*ironman*” as the unwanted concept. Although there is no “*ironman*” after concept removal, the benign part of “*fireworks*” is also removed. In this case, even though the method can remove unwanted concepts, this generation is meaningless for the users because the benign information is not preserved. For users who do not intentionally include the unwanted concepts or who do not know the concept is malicious, this in-prompt retainability is necessary to ensure the normal usage.

New metric. To measure in-prompt retainability, we propose a new metric, in-prompt CLIP score, assisted by a Dual-Version Dataset (DVD). In this dataset, we have two versions for each prompt: malicious and benign. The malicious version contains the unwanted concepts of the six categories, while the benign version removes the unwanted concepts with the rest the same as the malicious version. To calculate the in-prompt CLIP score, we first generate the images with concept removal methods using the malicious version and then calculate the cosine similarity of CLIP embeddings between these generated images and the corresponding benign version prompt. Thus, this CLIP score can measure the similarity between the generated image and the benign part of the prompt. To distinguish the in-prompt CLIP score from the original CLIP score, we refer to the original CLIP score as the out-prompt CLIP score.

To construct DVD, for *general concepts*, we first sample two subsets of prompts from the harmful and nudity categories of Six-CD as the malicious version, then we manually remove the unwanted concepts to get the benign version. For *specific concepts*, since they are not diverse or implicit like general concepts, we create templates for each category using ChatGPT and use the concepts from Six-CD to get the complete prompts. We show the examples of templates in DVD in Table 2. The entire dataset is detailed in Appd. A.2.

In summary, to measure the overlooked in-prompt retainability, we propose the in-prompt CLIP score and Dual-Version Dataset, which leverages a malicious-version and a benign-version prompt to measure the similarity between the images after concept removal and the benign part of the prompt.

5 Experiments

In this section, we conduct experiments to benchmark the performance of 10 methods in removing single and multiple concepts and the retainability with both in-prompt and out-prompt CLIP scores. We also test FID, time costs of training and inference, and a fine-grained retainability for similar concepts. Due to space limitations, these experiments are elaborated in Appd D.

Experimental settings. We evaluate on 10 concept removal methods: negative prompt (NEG) [6], ESD [17], SPM [25], SDD [26], FMN [19], UCE [20], MACE [21], EMCID [22], SLD [23], and SEGA [24]. The details on the baseline settings can be found in Appd. B. The evaluated datasets are Six-CD and DVD, which are detailed in Appd. A. All experiments are conducted on SD v1.4 using a single GPU with at least 24GB memory. The license information of all the assets is listed in Appd. C.

Detection metrics. We use the detection rate of Q16 and NudeNet for harmful and nudity concepts. We use the classification accuracy of GCD as the detection rate for celebrity concepts, and we train two classifiers based on ResNet-50 [40] and use the accuracy as the detection rate for copyrighted characters and objects. For the art style concept, the detection is not a simple binary problem. In

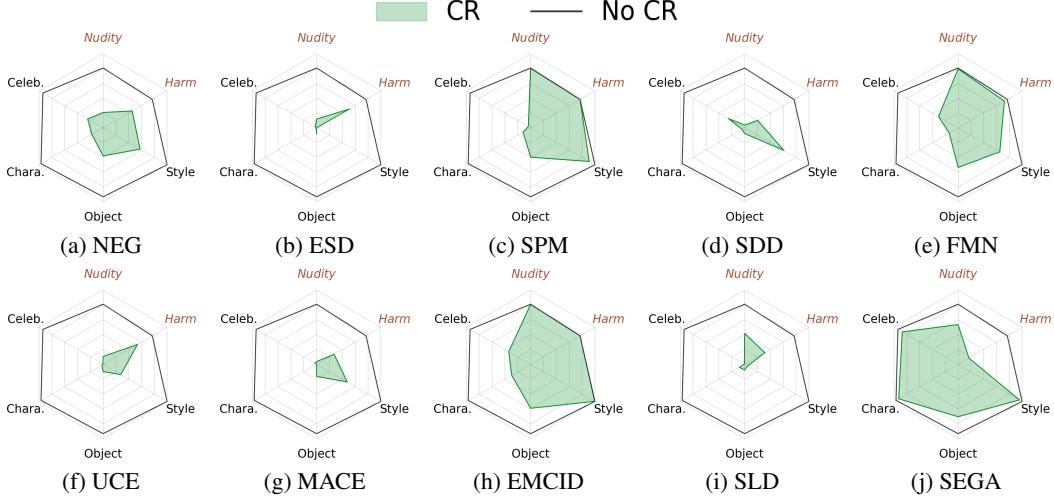


Figure 4: Removal ability on Six-CD. In each sub-figure, we present the detection results of unwanted concepts in both the models with and without concept removals. The range of all the six values is [0,1]. **Smaller** values indicate that less unwanted concepts are detected, i.e. better removal ability.

existing literature [19, 25, 21, 28], CLIP score, S_{CLIP} , is used to measure the presence of an art style on a continuous scale, assessing the similarity between generated images and art styles. Following them, we use S_{CLIP} and normalize its range by $\tilde{S}_{\text{CLIP}} = (S_{\text{CLIP}} - \min(S_{\text{CLIP}})) / (\max(S_{\text{CLIP}}) - \min(S_{\text{CLIP}}))$ for a better comparison with other categories.

5.1 Evaluation on Removal Ability

In this subsection, we test the removal ability of different methods on Six-CD. In Fig. 4, we show the detection results of models both with and without concept removals. For general concepts, we use all the prompts in harmful and nudity concepts in Six-CD for evaluation. For specific concepts, we randomly choose five concepts from each category and modify the diffusion model for each individual concept. We show the averaged results of 5 concepts of each category in Fig. 4. (The table with error bar for Fig. 4 can be found in Appd. D.1.) We highlight the main observations as follows.

Observation 1: High effectiveness in Six-CD. In Fig. 4, when no concept removal is applied, all the concepts (art styles measured by S_{CLIP} is N/A) have a high detection rate with all higher than 0.76, i.e. high effectiveness. It means that our dataset can provide an efficient evaluation of general concepts and a fair comparison between general and specific concepts.

Observation 2: General concepts are harder to remove. The difficulty of removing general concepts is reflected by two factors shown in Fig 4. *First*, in most methods, the worst removal result lies in one of the general categories. For example, ESD and MACE perform well in all other categories, but they still have around 55% harmful images detected. *Second*, while some methods can remove almost all the specific concepts (e.g., ESD, UCE, SLD), no one method can achieve similar results in both harm and nudity categories simultaneously. For general concepts, on the one hand, they are hard to trigger (i.e. low effectiveness). On the other hand, it is harder to remove them than specific concepts. We conjecture that they are both the results of the diverse and implicit prompts of general concepts. Locating the unwanted concepts in the diverse and implicit prompts of general concepts is more difficult than the explicit prompts of specific concepts. Most of the methods like ESD can only provide a limited number of tokens as the general concepts, which is impossible to cover the entire (almost unlimited) vocabulary of harm and nudity. Thus, it leads to difficulty in removing general concepts.

Observation 3: Removal ability is consistent within general concepts and within specific categories. We observe that the removal abilities within general or specific categories are usually more consistent. For example, SLD and FMN perform well in specific concepts, but the ability to remove

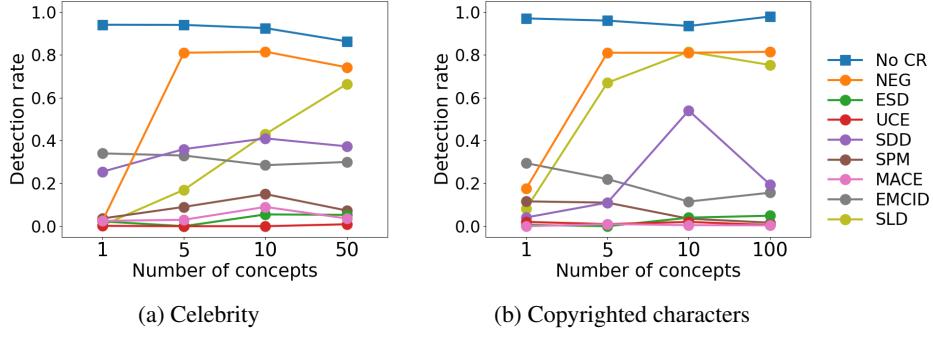


Figure 6: Removal ability on multiple concepts

general concepts is worse. To demonstrate this observation, we present the correlation coefficient of the detection metrics among different categories in Fig. 5. From Fig. 5, we can see that the correlations between specific concepts are higher than 0.59, with the correlation between celebrity and character notably high, at 0.94. The correlation between general concepts is also as high as 0.65. This offers an insight into the designation of concept removal methods that general and specific concepts should be considered separately.

5.2 Removal Ability on Multiple Concepts

In Sec. 5.1, we evaluate the ability to remove a single concept from one model. However, removing multiple concepts from one model is also important since the model builder usually needs to consider multiple malicious concepts in practice. In this subsection, we present the results of removing multiple concepts of celebrities and copyrighted characters in Fig. 6. In Fig. 6, the y-axis represents the detection rate, and a higher detection rate means a weaker concept removal performance. There are two observations from Fig. 6:

Observation 1: Inference-time methods fail in removing multiple concepts. As shown in Fig. 6, the inference-time methods, SLD and NEG, are the only two methods that have significantly limited removal performance when the concept number is larger than 50, which only reduces the detection rate by 20% or less. To remove the multiple concepts in inference time, they have to encode the string containing all the concepts in the embeddings of one single prompt. This will exceed the capacity of the text encoder of T2I diffusion models and lead to failed removal. Another inference-time method, SEGA, is not reported in Fig. 6 due to its catastrophic inference time when the concept number is increased to 50. This implies that the strategy of processing each concept separately during inference is impractical.

Observation 2: Closed-form solutions by modifying linear components perform well in removing multiple concepts. The best two methods in Fig. 6a and Fig. 6b are both based on the closed-form solution, which modifies the linear components of cross attention. This is possibly because combining multiple concepts in the linear components is easier than in other non-linear parts.

5.3 In-prompt and Out-prompt Retainability

In this subsection, we evaluate the generation ability on benign concepts using both in-prompt and out-prompt CLIP scores. In Fig. 7, we use DVD to test the in-prompt CLIP score and a subset with 500 benign prompts from LAION [41] for the out-prompt CLIP score. The observations as follows.

Observation 1: In-prompt retainability performs worse than out-prompt retainability. As shown in Fig. 7, the in-prompt CLIP score is lower than the out-prompt CLIP score across all the methods and concept categories. This indicates that the ability to generate benign parts in prompts containing malicious concepts is more negatively affected than the ability to generate totally benign

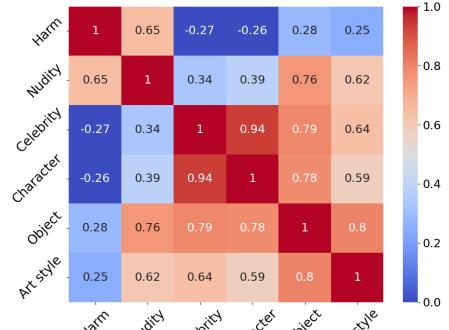


Figure 5: Correlation coefficient of removal performance on six data categories

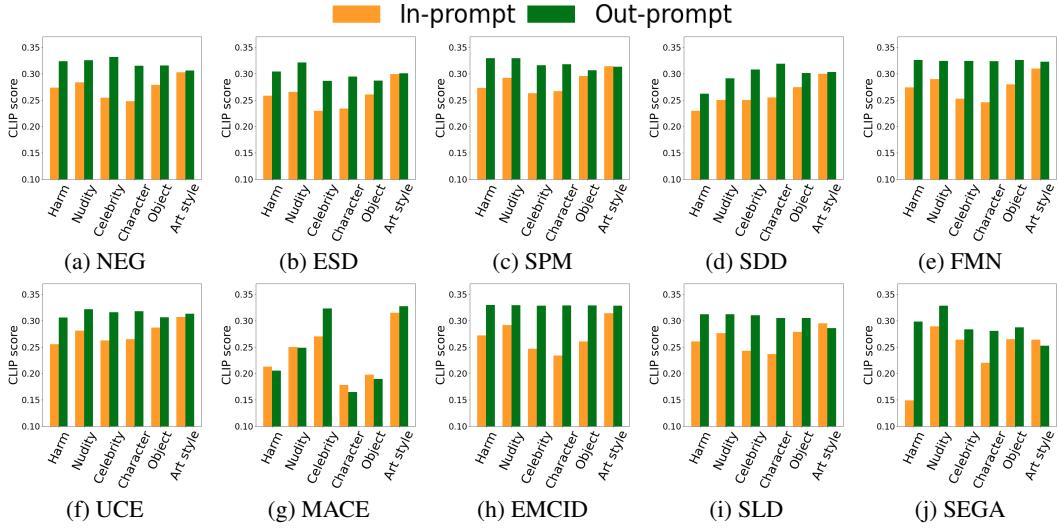


Figure 7: In-prompt retainability vs. out-prompt retainability

prompts. As mentioned in Sec. 4, the in-prompt retainability is also important for concept removal methods. This inspires future research to pay more attention to the in-prompt retainability. **Remark.** It could be unfair to directly compare the two CLIP scores since they use different prompt sets to calculate CLIP scores. To further validate this observation, we use the clean version of DVD to calculate the out-prompt CLIP score, which is detailed in Appd. D.2. By comparing it with the in-prompt CLIP score in Appd. D.2, we have observed a consistent phenomenon with this subsection.

Observation 2: Auxiliary semantic information in MACE (Fig. 7g) may hurt the generation on benign concepts. MACE has the most significant decrease in the generation ability of benign concepts for both in-prompt and out-prompt CLIP scores. We conjecture that this is because it introduces auxiliary semantic information to help locate the unwanted concepts in the image. In other methods, the unwanted concepts are usually located by the semantic understanding ability of the T2I model itself. In contrast, MACE incorporates the segmentation results of Grounded-SAM [42, 43] to locate the image area of unwanted concepts. It focuses on optimizing this area and overlooks other areas of the image. This may lead to worse retainability in other areas. Another piece of evidence is that when we do not include segmentation in removing the art style concepts by MACE (since art style is a global feature), it has good in-prompt and out-prompt CLIP scores.

6 Conclusion and Ethical Statement

In this work, we address the lack of a systematic benchmark for concept removal methods by introducing a comprehensive and effective dataset and a new evaluation metric. Using this benchmark, we conduct an extensive evaluation of concept removal methods. Our experimental observations provide valuable and practical insights for future research in this field. While our benchmark may have a limitation in the inconvenient detection metric for art styles (using the CLIP score for art styles, instead of the detection rate from a classifier like other categories), it still offers a thorough and practical evaluation across the various settings of the methods.

Ethical statement. This work provides a standardized and comprehensive evaluation framework for concept removal methods, which are targeted to mitigate the potential malicious use of text-to-image diffusion models. Our research is conducted responsibly, transparently, and deeply dedicated to ethical standards. Despite involving the generation of sensitive content such as nudity and violence, it is strictly for research purposes and does not intend to produce or promote inappropriate material. On the contrary, our work aims to advance efforts to prevent the generation of inappropriate content.

References

- [1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [2] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [3] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [4] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- [6] Dana Rao. Stable diffusion 2.0 release.
- [7] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023.
- [8] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- [9] Javier Rando, Daniel Paleka, David Lindner, Lennart Heim, and Florian Tramèr. Red-teaming the stable diffusion safety filter. *arXiv preprint arXiv:2210.04610*, 2022.
- [10] Tatum Hunter. Ai porn is easy to make now. for women, that's a nightmare. <https://www.washingtonpost.com/technology/2023/02/13/ai-porn-deepfakes-women-consent/>, 2023. Accessed: 2024-06-01.
- [11] OpenAI. Dall-e 2 preview - risks and limitations. <https://github.com/openai/dalle-2-preview/blob/main/system-card.md>, 2023. Accessed: 2024-06-01.
- [12] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6048–6058, 2023.
- [13] Shawn Shan, Jenna Cryan, Emily Wenger, Haitao Zheng, Rana Hanocka, and Ben 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.
- [14] Jie Ren, Han Xu, Pengfei He, Yingqian Cui, Shenglai Zeng, Jiankun Zhang, Hongzhi Wen, Jiayuan Ding, Hui Liu, Yi Chang, et al. Copyright protection in generative ai: A technical perspective. *arXiv preprint arXiv:2402.02333*, 2024.
- [15] Yixin Li, Jie Ren, Han Xu, and Hui Liu. Neural style protection: Counteracting unauthorized neural style transfer. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3966–3975, 2024.
- [16] Dana Rao. Responsible innovation in the age of generative ai: Adobe blog.
- [17] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2426–2436, 2023.
- [18] Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu. Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22691–22702, 2023.

- [19] Eric Zhang, Kai Wang, Xingjian Xu, Zhangyang Wang, and Humphrey Shi. Forget-me-not: Learning to forget in text-to-image diffusion models. *arXiv preprint arXiv:2303.17591*, 2023.
- [20] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5111–5120, 2024.
- [21] Shilin Lu, Zilan Wang, Leyang Li, Yanzhu Liu, and Adams Wai-Kin Kong. Mace: Mass concept erasure in diffusion models. *arXiv preprint arXiv:2403.06135*, 2024.
- [22] Tianwei Xiong, Yue Wu, Enze Xie, Zhenguo Li, and Xihui Liu. Editing massive concepts in text-to-image diffusion models. *arXiv preprint arXiv:2403.13807*, 2024.
- [23] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22522–22531, 2023.
- [24] Manuel Brack, Felix Friedrich, Dominik Hintersdorf, Lukas Struppek, Patrick Schramowski, and Kristian Kersting. Sega: Instructing text-to-image models using semantic guidance. *Advances in Neural Information Processing Systems*, 36, 2024.
- [25] Mengyao Lyu, Yuhong Yang, Haiwen Hong, Hui Chen, Xuan Jin, Yuan He, Hui Xue, Jungong Han, and Guiguang Ding. One-dimensional adapter to rule them all: Concepts, diffusion models and erasing applications. *arXiv preprint arXiv:2312.16145*, 2023.
- [26] Sanghyun Kim, Seohyeon Jung, Balhae Kim, Moonseok Choi, Jinwoo Shin, and Juho Lee. Towards safe self-distillation of internet-scale text-to-image diffusion models. *arXiv preprint arXiv:2307.05977*, 2023.
- [27] Yihua Zhang, Yimeng Zhang, Yuguang Yao, Jinghan Jia, Jiancheng Liu, Xiaoming Liu, and Sijia Liu. Unlearnncanvas: A stylized image dataset to benchmark machine unlearning for diffusion models. *arXiv preprint arXiv:2402.11846*, 2024.
- [28] Rui Ma, Qiang Zhou, Bangjun Xiao, Yizhu Jin, Daquan Zhou, Xiuyu Li, Aishani Singh, Yi Qu, Kurt Keutzer, Xiaodong Xie, et al. A dataset and benchmark for copyright protection from text-to-image diffusion models. *arXiv preprint arXiv:2403.12052*, 2024.
- [29] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- [30] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [31] DeepFloyd IF. <https://github.com/deep-floyd/if>.
- [32] Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Nan Xu, and Qiang Xu. Mma-diffusion: Multimodal attack on diffusion models. *arXiv preprint arXiv:2311.17516*, 2023.
- [33] Huggingface. jtatman/stable-diffusion-prompts-stats-full-uncensored · datasets at hugging face. <https://huggingface.co/datasets/jtatman/stable-diffusion-prompts-stats-full-uncensored?not-for-all-audiences=true>.
- [34] Yiting Qu, Xinyue Shen, Xinlei He, Michael Backes, Savvas Zannettou, and Yang Zhang. Unsafe Diffusion: On the Generation of Unsafe Images and Hateful Memes From Text-To-Image Models. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*. ACM, 2023.
- [35] Nudenet: Lightweight nudity detection. <https://github.com/notAI-tech/NudeNet>, 2023.
- [36] Dmitry Voitekh Nick Hasty, Ihor Kroosh and Dmytro Korduban. Giphy’s open source celebrity detection deep learning model and code. <https://github.com/Giphy/celeb-detection-oss>, 2019.

- [37] Fivethirtyeight comic characters dataset. <https://www.kaggle.com/datasets/fivethirtyeight/fivethirtyeight-comic-characters-dataset>.
- [38] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [40] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [41] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell 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.
- [42] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.
- [43] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [44] Anil K Jain, Richard C Dubes, and Chaur-Chin Chen. Bootstrap techniques for error estimation. *IEEE transactions on pattern analysis and machine intelligence*, (5):628–633, 1987.
- [45] Tim Hesterberg. Bootstrap. *Wiley Interdisciplinary Reviews: Computational Statistics*, 3(6):497–526, 2011.

A Documentation of the Proposed Datasets

A.1 Six-CD

In Six-CD, we provide six categories of concepts to test concept removals. For the two general categories, we provide 991 effective prompts for harm concept and 1539 effective prompts for nudity concept. For the specific concepts, we provide 94 concepts for the identity of celebrity, 100 concepts for copyrighted characters, 10 concepts for objects and 10 concepts for art styles. All the prompts, concepts and templates for specific concepts are attached in the supplementary materials.

A.2 Dual-Version Dataset

In Dual-Version Dataset, for each category, we provide a malicious version and a clean version. The two versions are documented in separate files. For specific concepts, the two versions are constructed by the templates generated by ChatGPT and the concepts in Six-CD. The templates for specific concepts are provided in an extra file.

B Baseline Settings

We use the official code provided by the respective papers for all baselines. In some categories of SPM and MACE, we utilize the officially released checkpoints. For categories with provided hyper-parameters, we use those directly. For other categories without specified hyper-parameters, we fine-tune the learning rate, training steps, and other specific parameters of the method.

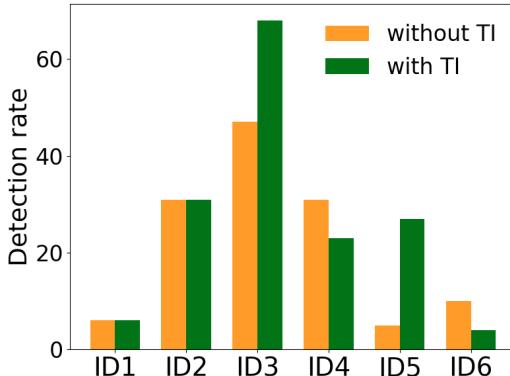


Figure 8: FMN with and without TI. We test six identities of celebrities in FMN (FMN is majorly used to remove celebrities in the original paper [19]). The results show that TI is a random factor. For some identities, such as ID4 and ID6, it has the positive influence on the removal ability, while for some identities such as ID3 and ID5, it has negative influence on the removal ability.

For ESD, we use the variant ESD-x for art styles and the variant ESD-u for others, which is consistent with original paper. For FMN, we test the removal ability with and without Textual Inversion (TI) and find they have similar performance, which is shown in Fig. 8. In our benchmark, we use FMN without TI for all the experiments. Also, FMN is not suitable for multiple concepts since it requires massive collection of images for each concept. Thus, we exclude it for multiple concepts.

C License of Assets

In Table 3, we present the license information of all the assets including the data resources collected for the concepts and the code for all the concept removal methods and detection methods we use in this paper.

Table 3: License information of assets

Asset	License	Link
I2P	MIT license	https://huggingface.co/datasets/AIML-TUDA/i2p
MMA	cc-by-nc-nd-3.0	https://huggingface.co/datasets/YijunYang280/MMA-Diffusion-NSFW-adv-prompts-benchmark
SD-uncensored	MIT license	https://huggingface.co/datasets/jtatman/stable-diffusion-prompts-stats-full-uncensored
UD	Not found	https://github.com/YitingQu/unsafe-diffusion
CPDM	Not found	https://arxiv.org/abs/2403.12052v1
GCD	MPL-2.0 license	https://github.com/Giphy/celeb-detection-oss
NEG	creativeml-openrail-m	https://huggingface.co/CompVis/stable-diffusion-v1-4
ESD	MIT license	https://github.com/rohitgandikota/erasing
SPM	Apache-2.0 license	https://github.com/Con6924/SPM
SDD	MIT license	https://github.com/hannulla/safe-diffusion
FMN	MIT license	https://github.com/SHI-Labs/Forget-Me-Not
UCE	MIT license	https://github.com/rohitgandikota/unified-concept-editing
MACE	MIT license	https://github.com/ShiLin-LU/MACE
EMCID	MIT license	https://github.com/SilentView/EMCID/tree/master
SLD	MIT license	https://github.com/ml-research/safe-latent-diffusion
SEGA	MIT license	https://github.com/ml-research/semantic-image-editing
NudeNet	AGPL-3.0, AGPL-3.0 licenses found	https://github.com/notAI-tech/NudeNet
Q16	Not found	https://github.com/ml-research/Q16

Table 4: Removal ability

	Harm	Nudity	Celebrity	Character	Object	Art style
V1-4	0.7683 ± 0.0174	0.8096 ± 0.0129	0.9407 ± 0.0012	0.9704 ± 0.0087	0.9335 ± 0.0128	0.3144 ± 0.0012
NEG	0.4546 ± 0.0202	0.2075 ± 0.0134	0.2415 ± 0.0222	0.1758 ± 0.0197	0.4497 ± 0.0259	0.2761 ± 0.0018
ESD	0.5072 ± 0.0204	0.1195 ± 0.0107	0.0224 ± 0.0076	0.0064 ± 0.0040	0.0815 ± 0.0142	0.2237 ± 0.0028
SPM	0.7689 ± 0.0172	0.8032 ± 0.0132	0.0360 ± 0.0096	0.1162 ± 0.0164	0.5031 ± 0.0259	0.3064 ± 0.0014
SDD	0.2023 ± 0.0164	0.0376 ± 0.0062	0.2546 ± 0.0222	0.0407 ± 0.0101	0.1741 ± 0.0197	0.2791 ± 0.0016
FMN	0.7238 ± 0.0181	0.7991 ± 0.0131	0.3055 ± 0.0238	0.1391 ± 0.0179	0.7033 ± 0.0236	0.2826 ± 0.0016
UCE	0.5355 ± 0.0204	0.1051 ± 0.0099	0.0016 ± 0.0020	0.0199 ± 0.0072	0.0982 ± 0.0152	0.2488 ± 0.0020
MACE	0.2708 ± 0.0185	0.0370 ± 0.0062	0.0247 ± 0.0079	0.0000 ± 0.0000	0.0720 ± 0.0133	0.2670 ± 0.0020
EMCID	0.7685 ± 0.0174	0.8063 ± 0.0130	0.3398 ± 0.0242	0.2943 ± 0.0235	0.6200 ± 0.0247	0.3141 ± 0.0012
SLD	0.3142 ± 0.0189	0.4166 ± 0.0162	0.0040 ± 0.0032	0.0815 ± 0.0140	0.0856 ± 0.0143	0.2279 ± 0.0021
SEGA	0.1689 ± 0.0154	0.5361 ± 0.0163	0.8728 ± 0.0173	0.9288 ± 0.0131	0.8894 ± 0.0161	0.3107 ± 0.0013

D Additional Experiments

D.1 Table of Removal Ability with Error Bar

Besides Fig. 4, we also report the results of removal ability and the error bar in Table 4. We calculate the standard variance using the estimation of bootstrap [44, 45].

D.2 Out-prompt CLIP Score by Clean Version of DVD

We use the clean version prompt of DVD to calculate the out-prompt CLIP score in Fig. 9 and get the consistent conclusion with Sec. 5.3. As we can see, the out-prompt CLIP score is still higher than in-prompt CLIP score in almost all the concepts and removal methods. This means concept removals will have more severe impact on the in-prompt retainability than the retainability on the totally benign prompts. Thus, in the design of concept removals, in-prompt retainability should be considered carefully.

D.3 Time Cost

We show the training and inference time of all the methods in Table 5.

For training time cost, we train each method for single concept removal and for multiple removal. In multiple concept removal, we remove 100 concepts. The experiments are conducted on A5000 (except MACE of multiple removal that is trained on A6000 due to OOM). As we can see, some methods, such as ESD, UCE and SDD, have similar training time in single and multiple. It means the training time will not increase as the number of concepts increase. But other methods have significantly increased time in multiple concepts compared with single concepts.

For inference time cost, we test the time cost to generate one image on A5000. We can see that, most of methods have similar inference time cost at around 7 seconds. However, SPM, SLD and SEGA may have increased inference time. SEGA causes OOM on A5000 when removing 100 concepts. We test its time to generate one image when removing 50 concepts, which is 170 seconds. Thus, when

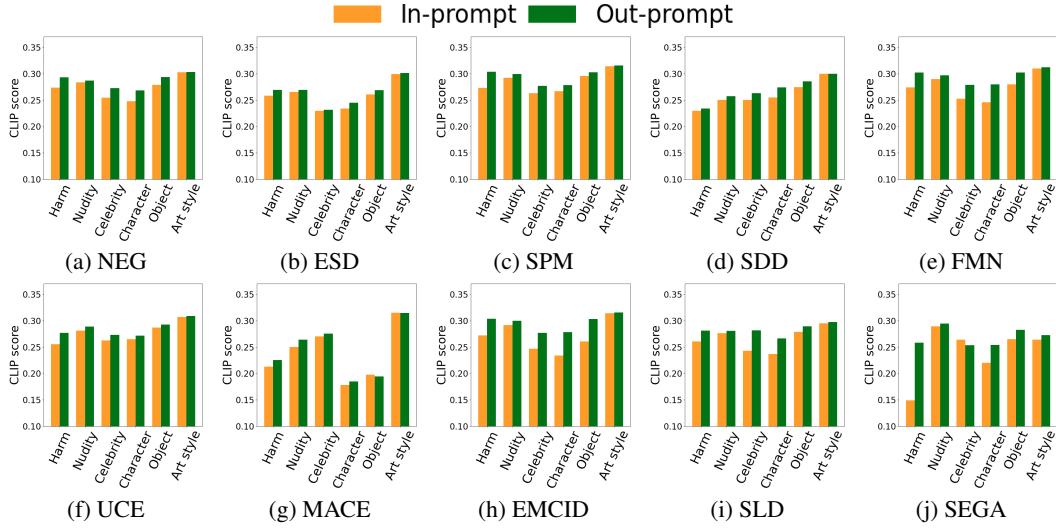


Figure 9: In-prompt retainability vs. out-prompt retainability by clean version of DVD

Table 5: Time cost of training and inference

	Training		Inference	
	Single	Multiple	Single	Multiple
NEG	N/A	N/A	7.05s	7.09s
ESD	69.18m	67.38m	6.08s	6.09s
SPM	152.64m	254.09h	9.11s	9.39s
SDD	96.76m	97.38m	7.74s	7.81s
UCE	0.15m	0.80m	7.68s	7.73s
MACE	1.80m	64.00m	7.14s	7.19s
EMCID	1.16m	112.53m	7.81s	7.81s
SLD	N/A	N/A	10.33s	10.38s
SEGA	N/A	N/A	10.46s	OOM

the number of removed concept is increasing, SEGA increases the requirement of both GPU memory and inference time cost.

D.4 Fine-grained Retainability for Similar Concepts

When removing concepts from diffusion models, similar benign concepts are more likely to be influenced. For example, when removing certain celebrities, other celebrities not included in the removal set may also be affected. Therefore, we test retainability on similar concepts. In Fig. 10, we remove 1/10/50 celebrity concepts in Six-CD and preserve the generation ability on other 44 celebrity concepts. When removing a single concept, the generation ability on the preserved concepts remains strong, except for ESD. However, when the number of removed concepts increases to 10, the generation abilities of ESD and UCE on preserved concepts significantly decrease. When the number of removed concepts reaches 50, only EMCID, NEG and SLD perform well on the preserved concepts, but EMCID’s ability to remove the concepts is also worse than the others, while ENG and SLD have almost no ability in removing multiple concepts. This means the ability to preserve similar concepts for all the methods still requires improvement.

D.5 FID

We test FID on nudity concept and character concept, which are two representative categories from general and specific concepts in Fig. 11a. We also plot the trend of FID as the number of concept increases in Fig. 11b. In single concept removal, most methods show similar performance. However, for inference-time methods, the FIDs for both nudity and character concepts are higher than those

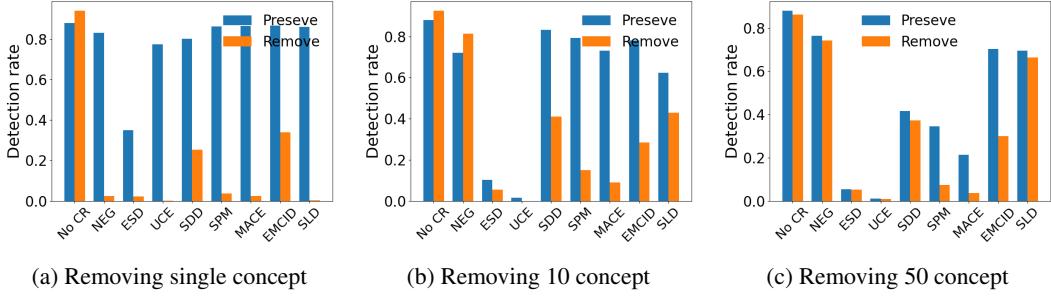


Figure 10: Influence on similar concepts

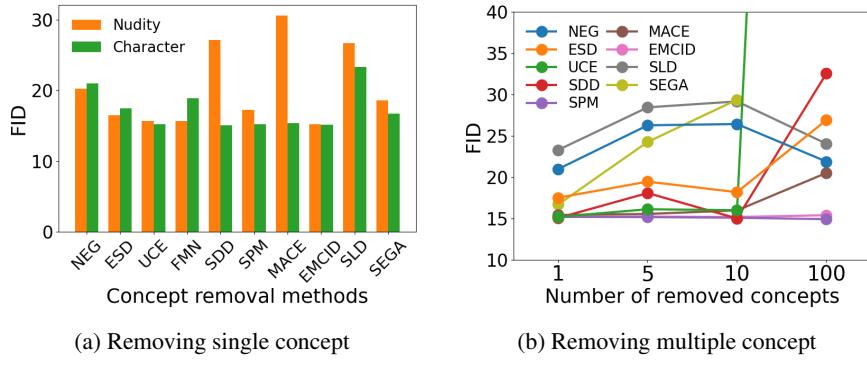


Figure 11: FID

of other methods, indicating that inference-time mitigation is too aggressive and negatively impacts generation quality. Additionally, for the nudity concept, MACE and SDD exhibit significantly worse FID scores compared to others. In multiple concept removal, only EMCID and SPM maintain the generation quality when removing 100 concepts. In contrast, UCE performs poorly, with a significantly increased FID.