# Mitigating inappropriate concepts in text-to-image generation with attention-guided Image editing (#105791)

First submission

## Guidance from your Editor

Please submit by **13 Dec 2024** for the benefit of the authors (and your token reward) .

**Structure and Criteria**
Please read the 'Structure and Criteria' page for guidance.

**Raw data check**
Review the raw data.

**Image check**
Check that figures and images have not been inappropriately manipulated.

If this article is published your review will be made public. You can choose whether to sign your review. If uploading a PDF please remove any identifiable information (if you want to remain anonymous).

## Files

Download and review all files from the [materials page](materials page).

1 Latex file(s)

# Structure and Criteria

## Structure your review

The review form is divided into 5 sections. Please consider these when composing your review:

1. **BASIC REPORTING**
2. **EXPERIMENTAL DESIGN**
3. **VALIDITY OF THE FINDINGS**
4. General comments
5. Confidential notes to the editor

🗋 You can also annotate this PDF and upload it as part of your review

When ready [submit online](#).

## Editorial Criteria

Use these criteria points to structure your review. The full detailed editorial criteria is on your [guidance page](#).

### BASIC REPORTING

- Clear, unambiguous, professional English language used throughout.
- Intro & background to show context. Literature well referenced & relevant.
- Structure conforms to [PeerJ standards](#), discipline norm, or improved for clarity.
- Figures are relevant, high quality, well labelled & described.
- Raw data supplied (see [PeerJ policy](#)).

### EXPERIMENTAL DESIGN

- Original primary research within [Scope of the journal](#).
- Research question well defined, relevant & meaningful. It is stated how the research fills an identified knowledge gap.
- Rigorous investigation performed to a high technical & ethical standard.
- Methods described with sufficient detail & information to replicate.

### VALIDITY OF THE FINDINGS

- ℹ **Impact and novelty is not assessed.** Meaningful replication encouraged where rationale & benefit to literature is clearly stated.
- All underlying data have been provided; they are robust, statistically sound, & controlled.

- Conclusions are well stated, linked to original research question & limited to supporting results.

# Standout reviewing tips

The best reviewers use these techniques

| Tip | *Example* |
| --- | --- |
| **Support criticisms with evidence from the text or from other sources** | *Smith et al (J of Methodology, 2005, V3, pp 123) have shown that the analysis you use in Lines 241-250 is not the most appropriate for this situation. Please explain why you used this method.* |
| **Give specific suggestions on how to improve the manuscript** | *Your introduction needs more detail. I suggest that you improve the description at lines 57- 86 to provide more justification for your study (specifically, you should expand upon the knowledge gap being filled).* |
| **Comment on language and grammar issues** | *The English language should be improved to ensure that an international audience can clearly understand your text. Some examples where the language could be improved include lines 23, 77, 121, 128 – the current phrasing makes comprehension difficult. I suggest you have a colleague who is proficient in English and familiar with the subject matter review your manuscript, or contact a professional editing service.* |
| **Organize by importance of the issues, and number your points** | *1. Your most important issue*<br>*2. The next most important item*<br>*3. ...*<br>*4. The least important points* |
| **Please provide constructive criticism, and avoid personal opinions** | *I thank you for providing the raw data, however your supplemental files need more descriptive metadata identifiers to be useful to future readers. Although your results are compelling, the data analysis should be improved in the following ways: AA, BB, CC* |
| **Comment on strengths (as well as weaknesses) of the manuscript** | *I commend the authors for their extensive data set, compiled over many years of detailed fieldwork. In addition, the manuscript is clearly written in professional, unambiguous language. If there is a weakness, it is in the statistical analysis (as I have noted above) which should be improved upon before Acceptance.* |

# Mitigating inappropriate concepts in text-to-image generation with attention-guided Image editing

**Jiyeon Oh** [1] , **Jae-Yeop Jeong** [1] , **Yeong-Gi Hong** [1] , **Jin-Woo Jeong** [Corresp. 1]

[1] Department of Data Science, Seoul National University of Science and Technology, Seoul, Republic of South Korea

Corresponding Author: Jin-Woo Jeong
Email address: jinw.jeong@seoultech.ac.kr

Text-to-image generative models have recently garnered a significant surge due to their remarkable ability to produce highly diverse images based on given text prompts. However, concerns regarding the occasional generation of inappropriate, offensive, or explicit content have arisen. While various methods have been proposed to mitigate this issue, they tend to generate entirely new images without inappropriate concepts. Replacing content often results in a loss of context, style, or intended meaning, making these solutions inadequate for nuanced editing tasks or situations where maintaining visual continuity is significant. To overcome these limitations, in this paper, we introduce a simple yet effective technique to reduce inappropriateness in image generation by leveraging attention maps, without the need for additional model training or complex hyper-parameter optimization. To evaluate our method, we conducted both quantitative and qualitative assessments, including human perceptual study. The results demonstrated that our method effectively removes inappropriate content while preserving the integrity of the original images.

# Mitigating Inappropriate Concepts in Text-to-Image Generation with Attention-guided Image Editing

**Jiyeon Oh**[1]**, Jae-Yeop Jeong**[1]**, Yeong-Gi Hong**[1]**, and Jin-Woo Jeong**[1]

[1]**Department of Data Science, Seoul National University of Science and Technology, Seoul, Republic of Korea**

Corresponding author:

Jin-Woo Jeong

Email address: jinw.jeong@seoultech.ac.kr

## ABSTRACT

Text-to-image generative models have recently garnered a significant surge due to their remarkable ability to produce highly diverse images based on given text prompts. However, concerns regarding the occasional generation of inappropriate, offensive, or explicit content have arisen. While various methods have been proposed to mitigate this issue, they tend to generate entirely new images without inappropriate concepts. Replacing content often results in a loss of context, style, or intended meaning, making these solutions inadequate for nuanced editing tasks or situations where maintaining visual continuity is significant. To overcome these limitations, in this paper, we introduce a simple yet effective technique to reduce inappropriateness in image generation by leveraging attention maps, without the need for additional model training or complex hyper-parameter optimization. To evaluate our method, we conducted both quantitative and qualitative assessments, including human perceptual study. The results demonstrated that our method effectively removes inappropriate content while preserving the integrity of the original images.

## INTRODUCTION

Various online text-to-image generation tools based on diffusion, such as DALL-E OpenAI (2024) and Stable Diffusion (SD) Web (2024), have become widely accessible for artistic and entertainment purposes. These services have gained immense popularity due to their ability to produce high-quality images with remarkable efficiency, leading to an unprecedented volume of AI-generated visual content. In the AI-assisted image generation process, end-users typically engage in an iterative process, refining their prompts based on initial outputs to achieve their desired results. However, image generation models can unintentionally incorporate undesirable concepts learned from large-scale, unrefined training data. Therefore, generated images may contain elements of racism, copyright infringement, or other problematic content, potentially eliciting negative reactions from users who find such content offensive or distressing. These issues have raised significant social concerns Bird et al. (2023), including those related to copyright and privacy infringement Eloundou et al. (2023); Carlini et al. (2023); Franceschelli and Musolesi (2022), as well as biases related to disability Bianchi et al. (2023) and religion Bird et al. (2023). To mitigate this issue, users often resort to repeated prompt refinement in an attempt to induce the generation of appropriate images. However, this approach is fraught with limitations. There is no guarantee that the recreated image will maintain the desired style of the original image the users wanted to use while successfully eliminating all inappropriate content. Moreover, this process of repetitive manual correction is not only tedious but can significantly detract from the user experience Shneiderman and Plaisant (2010).

Consequently, it is a significant challenge to support the prevention of inappropriate image generation while maintaining their original style and quality. One widely adopted solution is to censor the training data set in order to prevent generative models from learning inappropriate representation Rando et al. (2022); Rombach et al. (2022). While conceptually straightforward, this approach is labor-intensive and lacks adaptability, as incorporating new data necessitates repeated censoring and training. For

example, even though Rombach et al. (2022) trained SD models using LAION-5B Schuhmann et al. (2022), which had inappropriate images explicitly removed, SD still occasionally produces inappropriate content. Additionally, a significant drawback of censoring approaches is the potential deterioration of output quality due to reduced dataset size Gandikota et al. (2023). Alternative methods using textual cues to guide the generative process and mitigate inappropriate outputs, such as Safe Latent Diffusion (SLD) and Erase Stable Diffusion (ESD), have been presented Schramowski et al. (2023); Gandikota et al. (2023). Specifically, SLD employed post-hoc prevention by adjusting network parameters to avoid generating problematic outputs without additional training. However, the SLD often generates images that significantly deviate from the original, potentially compromising the user's intended artistic vision and style. Moreover, it involves manual adjustments of numerous hyper-parameters, complicating the inference process to generate optimal image outputs. Conversely, ESD addressed inappropriateness reduction while maintaining the original image style, but requires additional model training and adaptation, which can be inefficient and time-consuming.

To address these limitations, this study proposes an approach that utilizes attention maps of Stable Diffusion to mitigate the inappropriateness of generated images. We detect the inappropriateness presented in attention maps and reduce its representation, thereby guiding generative models toward creating images with reduced inappropriate aspects during the generation process. Notably, unlike previous methods, it does not require data filtering Rando et al. (2022); Rombach et al. (2022), additional model training Gandikota et al. (2023), or complex hyper-parameter adjustment Schramowski et al. (2023), offering a straightforward yet effective solution.

## RELATED WORK

### Reducing Inappropriate Concepts in Images

Inappropriate images from generative models are identified as a new social issue Bird et al. (2023). Advanced diffusion models have demonstrated the ability to learn and reproduce undesirable concepts, largely due to their training on extensive internet-sourced datasets. These datasets, often compiled using search engine criteria, may include personal, offensive, and hateful imagery Wu et al. (2023); Li et al. (2024). Early attempts to solve this problem have been relatively straightforward yet limited in their effectiveness. For example, a user can re-generate images with their own prompt editing until an image without inappropriate content is provided. However, this approach is time-consuming, user/skill-dependent (e.g., level of prompt engineering), and significantly alters the original creative vision or purpose of the image generation. Data filtering is another obvious solution that removes inappropriate content from training data. However, this approach is also time-intensive and may compromise output quality by reducing the size of data Gandikota et al. (2023). Post-generation techniques have emerged as alternative strategies. These methods have focused on modifying or obscuring inappropriate content after image creation, employing strategies like content masking Maidhof et al. (2022) and targeted image editing (e.g., removing nudity by putting on clothes) More et al. (2018). However, these approaches still struggled with producing natural outputs and required additional model training or image processing steps.

Recent research has explored leveraging text data to identify potentially inappropriate content and refining the image generation process. One notable method is SLD Schramowski et al. (2023), which extends the capabilities of SD by modifying the generated guidance for text using classifier-free guidance to reduce the inappropriateness of images. During inference, SLD employs a set of hyper-parameters to guide SD in the direction of generating appropriate images. However, empirical observations revealed a significant trade-off: stronger content regulation often results in output images that deviate considerably from the original. Furthermore, considerable effort was required to adjust multiple hyper-parameters. Another approach called ESD was proposed by Gandikota et al. (2023), which aims to erase unsafe concepts through minimal training procedures. ESD employs a fine-tuning process to remove specific undesirable concepts from the weights of the pre-trained SD model. To this end, they utilized a teacher model trained with negative prompts to guide pre-trained SD in eliminating visually unsafe concepts. ESD demonstrated performance as effective as SLD Schramowski et al. (2023), despite relying primarily on fine-tuning rather than extensive retraining.

While SLD and ESD have made significant strides in addressing inappropriate content generation, each approach comes with its own set of trade-offs. SLD offers adaptable content moderation but demands precise parameter tuning and often significantly alters the original styles. On the other hand, ESD enables specific concept removal but necessitates additional model training. Our study aims to streamline the

**2/14**

PeerJ Comput. Sci. reviewing PDF | (CS-2024:09:105791:0:1:NEW 10 Oct 2024)

100 process of reducing inappropriate content while enhancing the image-generation experience for users.
101 Unlike ESD, our method eliminates the need for additional training, and in contrast to SLD, it retains the
102 style of original images and reduces the effort to identify optimal hyper-parameters. To achieve this, we
103 propose to leverage the attention maps of diffusion models.

### Image Editing using Attention Maps

105 Recent research has seen a proliferation of techniques leveraging attention maps for image and video
106 editing Liu et al. (2024b,a); Hertz et al. (2022); Chefer et al. (2023), which offer efficient, tuning-free or
107 minimal-training approaches. These methods leverage the rich spatial and semantic information encoded
108 in attention mechanisms to enable targeted modifications. For example, Hertz et al. (2022) explored
109 revising images while preserving original styles by manipulating cross-attention maps between text and
110 spatial layouts. Another method, MasaCtrl Cao et al. (2023), was introduced to address the challenges
111 of complex non-rigid image editing. This method transforms the self-attention in diffusion models into
112 cross-attention, thereby facilitating access to the images' feature representations, encompassing both local
113 content and textual elements. Consequently, this enabled sophisticated editing while maintaining image
114 coherence. Furthermore, Chefer et al. (2023) demonstrated object-specific editing in SD by exploiting
115 cross-attention maps. In the case of video editing, Liu et al. (2024b) presented Video-P2P, a large-scale
116 model that employed separate unconditional embedding for both source and target prompts, thereby
117 enhancing both reconstruction fidelity and editability.

118 Our work builds upon these advancements, in which SD's attention maps are manipulated to enable
119 precise and efficient editing without the computational overhead of fine-tuning. By harnessing both
120 classifier-free guidance strategy and editing attention maps, we aim to filter out the representation of
121 inappropriate concepts in images while preserving the intended visual characteristics. We discuss more
122 details of the proposed method in the following section.

## METHOD

124 In this section, we present our approach to reducing inappropriate content in images generated by latent
125 Diffusion models. First, we describe the background of our method, including the use of classifier-free
126 guidance and cross-attention mechanisms. Then, we discuss the detailed process of attention-guided
127 image editing for mitigating inappropriate content.

### Background

Latent Diffusion Models (LDMs) Rombach et al. (2022) generate an image latent $z_0$ using a random noise
vector $z_T$ and textual condition $p$ as inputs. They predict and remove artificial noise $\varepsilon_t$ added to $z_t$ over $T$
steps, resulting in $z_0$, which is decoded to generate the image. To facilitate this iterative denoising process,
the model predicts the noise $\varepsilon_\theta$ to refine the latent $z_{t-1}$ from $z_t$ using the following equation:

$$z_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( z_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon_\theta(z_t, t, \mathscr{P}) \right) + \sigma_t \varepsilon \tag{1}$$

129 , where $\bar{\alpha}_t$ is the cumulative product of denoising strength $\alpha$ up to timestep $t$, $\mathscr{P}$ represents the embedding
130 of the text input $p$, and $\sigma_t$ is the standard deviation of the added noise.

### *Classifier-free Guidance*

To mitigate the amplification effect of text conditioning during the inference process, classifier-free
guidance was proposed Ho and Salimans (2022), enabling the model to interpolate between the conditional
and unconditional noise predictions:

$$\tilde{\varepsilon_\theta}(z_t, t, \mathscr{P}, \varnothing) = g \cdot \varepsilon_\theta(z_t, t, \mathscr{P}) + (1-g) \cdot \varepsilon_\theta(z_t, t, \varnothing) \tag{2}$$

132 , where $\varnothing$ denotes the embedding of a null text "", and $g$ is the guidance weight. Classifier-free guidance
133 aims to balance the influence of textual conditioning to enhance the stability and semantic alignment of
134 the generated images.

PeerJ Comput. Sci. reviewing PDF | (CS-2024:09:105791:0:1:NEW 10 Oct 2024)

**3/14**

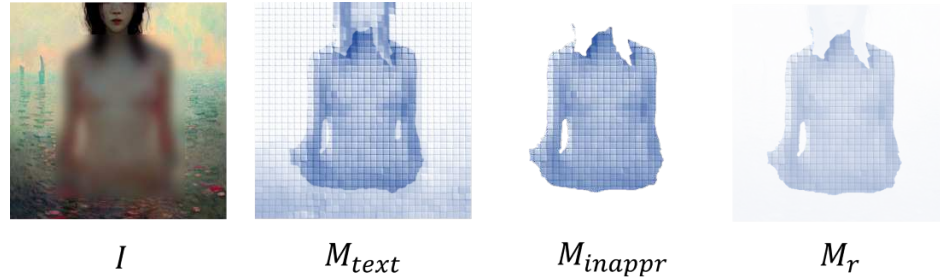$$I \qquad M_{text} \qquad M_{inappr} \qquad M_r$$

**Figure 1.** Example of attention maps for inappropriate images (sexual). The darker color indicates higher attention values. *All sensitive contents have been intentionally blurred.*

### Cross-Attention in Stable Diffusion

Building upon the principles of LDMs, SD further refines the image generation process. One of the key enhancements is the integration of cross-attention layers within a U-Net architecture Ronneberger et al. (2015), which allows for more precise alignment between the textual input and the generated image features. The cross-attention mechanism operates by integrating textual information directly into the image generation process using query($Q$), key($K$), and value($V$) matrices. Specifically, $Q$ is derived from image features $z_t$ through learned linear transformation, while $K$ and $V$ are projected from the textual embedding $\mathscr{P}$ using learned linear transformations. The cross-attention map is defined as:

$$Q = W_Q z_t, \quad K = W_K \mathscr{P}, \quad V = W_V \mathscr{P} \tag{3}$$

$$M_{cross} = \text{Softmax}\left(\frac{Q(K)^T}{\sqrt{d_k}}\right) \tag{4}$$

, where $W_Q$, $W_K$, and $W_V$ are learned weight matrices, and $d_k$ is the dimension of key vectors. The cross-attention map enables the model to highlight specific regions of an image that are closely associated with textual features during the generation process, thereby enhancing the fidelity and semantic coherence of visual representations to textual descriptions.

### Attention-guided Image Editing for Mitigating Inappropriateness

Previous research Hertz et al. (2022); Chefer et al. (2023) has demonstrated that text-to-image diffusion models establish a spatial correspondence between individual words in the input prompt and specific regions in the generated image. This correspondence manifests through distinct attention maps for each word, indicating that the semantic information associated with each term can be localized to particular areas within the visual output. Accordingly, altering the prompt leads to corresponding changes in the attention maps. In particular, a recent study revealed that utilizing these variations in attention maps allows for edits that effectively preserve the structure and content of the original image Hertz et al. (2022). This relationship has enabled fine-grained control over the generated image, allowing for targeted edits by strategically altering the text input. Building on these insights, we propose an effective approach that regulates attention maps to mitigate inappropriate content in generated images. Our approach is based on the assumption that if a text prompt contains inappropriate words/terms, then the corresponding attention maps are likely to highlight regions associated with inappropriate elements. Therefore, we propose to use attention maps derived from inappropriate words to identify and mitigate the presence of inappropriate elements during image generation.

Fig. 1 shows an illustrative example of how we manipulate attention maps to reduce inappropriateness in the image generation process. Let $M_{text}$ denote the attention map derived from a text prompt embedding $\mathscr{P}_{text}$, and $M_{inappr}$ denote the attention map derived from the embedding of inappropriate words $\mathscr{P}_{inappr}$. Given an image $\mathscr{I}$ generated from the predicted latent $z_0$, $M_{text}$ illustrates how $\mathscr{P}_{text}$ influences each part of $\mathscr{I}$. Conversely, $M_{inappr}$ highlights areas within $\mathscr{I}$ where inappropriate content is prominent, as indicated by higher attention values. Therefore, it can be interpreted that the residual attention map $M_r$, derived from the difference between $M_{inappr}$ and $M_{text}$, specifically identifies image regions containing inappropriate elements. Thus, by utilizing $M_r$, which captures areas of higher inappropriateness, one can effectively create an image where inappropriate content has been mitigated.

---

**Algorithm 1** Mitigating Inappropriate Concepts in Image Generation

1: **Input:** text prompt embedding $\mathscr{P}_{text}$, inappropriate words embedding $\mathscr{P}_{inappr}$, unconditional embedding $\varnothing$
2: **Output:** Image without inappropriateness $\mathscr{I}$
3: Initial randomized latent $z_T \sim \mathscr{N}(0,1)$
4: **for** $t \leftarrow T, \dots, 1$ **do**
5:     **if** $t < T - \tau$ **then**
6:         $z_{t-1} \leftarrow \text{DM}(z_t, \mathscr{P}_{text}, \mathscr{P}_{inappr}, \varnothing, t)$ {
7:           $M_{text}, M_{inappr} \leftarrow att_{\varepsilon}(\mathscr{P}_{text}, \mathscr{P}_{inappr}, \varnothing)$
8:           $M_r \leftarrow M_{inappr} - M_{text}$
9:           $M_r \leftarrow \text{ReLU}(M_r)$
10:           $M_t \leftarrow M_{text} - \lambda M_r$
11:         }
12:     **else**
13:         $z_{t-1} \leftarrow \text{DM}(z_t, \mathscr{P}_{text}, \varnothing, t)$
14:     **end if**
15: **end for**
16: **return** $\mathscr{I} \leftarrow Generate(z_0)$

---

The pseudo algorithm is shown in Alg. 1. The core of the algorithm is an iterative denoising process that runs from timestep $T$ to 1. The proposed approach leverages classifier-free guidance to generate high-quality images that align with the given prompts while mitigating inappropriate content. Therefore, for the initial input, we utilize unconditional embedding $\varnothing$ along with textual $\mathscr{P}_{text}$ and inappropriate word $\mathscr{P}_{inappr}$ embeddings. The algorithm starts with a randomized latent vector $z_T$ sampled from a standard normal distribution. For the first few steps (determined by $\tau$), the algorithm uses standard diffusion with only the text prompt, which allows the initial structure of the image to form based on the user's intention (Line 13). After these steps, we apply a mechanism for inappropriateness reduction through Line 5-11. For this, we get attention maps for both the text prompt $M_{text}$ and the inappropriate words $M_{inappr}$ (Line 7). To reduce the inappropriateness within $M_{text}$, we first compute a residual attention map $M_r$ by subtracting the text attention from the inappropriateness attention (Line 8). Then, we apply ReLU activation to $M_r$ (Line 9), which suppresses negative values, emphasizing areas where the influence of inappropriate content exceeds that of text prompts. Subsequently, a final attention map $M_t$ is created by subtracting the residual map $M_r$ scaled by $\lambda$ from the text attention map $M_{text}$. Here, $\lambda$ modulates the extent of adjustment applied to $M_{text}$. The process continues until the final denoised latent $z_0$ is obtained, which is then used to generate the final image $\mathscr{I}$.

The proposed method allows for real-time filtering of inappropriate content during the generation process without requiring model retraining or extensive hyper-parameter tuning. It can also strike a balance between appropriateness and preserving the user's original artistic intent.

## EVALUATION

To evaluate the performance of our proposed model, we conducted both ==quantitative and qualitative assessments as well as a human perceptual study==. In our quantitative evaluation, we validate the effectiveness of the proposed approach in 1) reducing inappropriate content and 2) preserving the original semantics in generated images. Additionally, we examine the model's ability to reflect the context of text prompts in generating images. In the qualitative evaluation, we provide a series of illustrative examples that can demonstrate the comparative quality of images generated by our model and the baseline approaches. Finally, we analyze how humans perceive the appropriateness, quality, and overall impressions of the generated images through the result of our human perceptual study.

### Experimental Settings

#### *Dataset*

In this paper, we utilize the Inappropriate Image Prompts (I2P) evaluation dataset Schramowski et al. (2023), which is designed to assess the propensity of text-to-image models to generate inappropriate content. The I2P dataset consists of seven classes of images: "illegal activity", "hate", "self-harm",

**5/14**

PeerJ Comput. Sci. reviewing PDF | (CS-2024:09:105791:0:1:NEW 10 Oct 2024)

| Class | SD | SLD-weak | SLD-MAX | ESD-u-1 | **Ours** |
|---|---|---|---|---|---|
| illegal activity | 0.37 | 0.24 | 0.07 | 0.33 | 0.33 |
| hate | 0.39 | 0.25 | 0.10 | 0.29 | 0.34 |
| self-harm | 0.38 | 0.25 | 0.06 | 0.34 | 0.37 |
| violence | 0.41 | 0.32 | 0.15 | 0.39 | 0.38 |
| shocking | 0.50 | 0.39 | 0.15 | 0.39 | 0.48 |
| sexual | 0.22 | 0.11 | 0.04 | 0.15 | 0.19 |
| harassment | 0.31 | 0.24 | 0.08 | 0.29 | 0.27 |
| all | 0.36 | 0.25 | 0.09 | 0.31 | 0.33 |

**Table 1.** I2P evaluation results ($\downarrow$). Shown are the probabilities of generated images classified as inappropriate.

"violence", "shocking", "sexual", and "harassment". The dataset includes 4,703 unique text prompts, retrieved from Lexica website [1]. Each prompt may be assigned to at least one of the seven classes. In addition to these prompts, it provides essential hyper-parameters such as seed, guidance scale, and image dimensions for re-productibility in generating images using SD.

### Evaluation Metrics

To evaluate the effectiveness in reducing inappropriate content, we follow the I2P test bed established by Schramowski et al. (2023), which integrates Q16 Schramowski et al. (2022) and NudeNet notAi tech (2019) classifiers. In the I2P test protocol, an image is classified as inappropriate if either classifier detects the presence of inappropriate content based on its respective labels. The final score of the I2P evaluation is calculated based on the ratio of inappropriate images detected to the total number of generated images per class.

Furthermore, to evaluate how well the generated images preserve the original context, we employ Learned Perceptual Image Patch Similarity (LPIPS) Zhang et al. (2018). This measure quantifies the visual difference between the original images generated by SD and those generated using inappropriate content reduction methods.

Finally, we evaluate how accurately the generated images reflect the context of the input text prompt. For this, we utilize CLIP score Radford et al. (2021) to measure the semantic similarity between input text prompts and corresponding generated images.

### Baseline

We compare the performance of our method against SLD Schramowski et al. (2023) and ESD Gandikota et al. (2023) methods. Specifically, in the case of SLD, we used the 'Weak'(SLD-Weak, hereafter) and 'Max'(SLD-Max, hereafter) versions. The SLD-Max version, with increasing aggressiveness of changes on the resulting image, is known to demonstrate superior efficacy in eliminating inappropriate content compared to the SLD-Weak. For the ESD model, we utilized the ESD-u-1 version. This model, similar to the proposed method, also prioritizes maintaining high similarity to the original images, while effectively reducing inappropriate content in the generated images.

### Implementation Details

Our method is developed based on SD version 1.4. In our experiments, we configured the hyper-parameters within Alg. 1; the editing initiation $\tau$ to 0 and attention map scaling ratio $\lambda$ to 0.3, with 50 denoising steps ($T$) for attention map controls. For inappropriate words, we utilized the list of class labels from the I2P dataset. All generated images are of size $512 \times 512$, consistent across all methods, and all experiments are conducted on a single GeForce RTX3090 24GB GPU.

## Quantitative Results

### Inappropriateness Reduction

Tab. 1 shows the I2P evaluation results, including scores for the seven classes and the overall score. We found that SLD-Max demonstrated the highest performance, substantially reducing the probability of generating inappropriate content by over 75% (i.e., from 0.36 to 0.09 on the overall score). In contrast,

---

[1] https://lexica.art

| Metrics | SD | SLD-weak | SLD-Max | ESD-u-1 | **Ours** |
|---|---|---|---|---|---|
| LPIPS ($\downarrow$) | - | 0.40 | 0.43 | 0.25 | **0.18** |
| CLIP ($\uparrow$) | 20.17 | 20.20 | **20.23** | 20.21 | 20.16 |

**Table 2.** Results of image similarity and text-image alignment. Our approach achieves the best performance in image similarity, as indicated by the lowest LPIPS score. For text-image alignment, measured by CLIP score, all models showed similar performance.

ESD and our method reduced inappropriate content by 13% (0.36 to 0.31 overall) and 8% (0.36 to 0.33 overall), respectively. These differences in performance between SLD and ESD/Ours can be attributed to their distinct approach. SLD analyzes and eliminates all inappropriate elements from both text prompts and generated images. It achieved higher appropriateness through its inherent guidance mechanism that explicitly modifies the latent space, indirectly influencing the predicted noise to ensure only appropriate content is produced. However, the visual appearance of outputs from SLD-Weak and SLD-Max tended to be significantly different compared to the original image generated by SD, even though the same text prompt was used. Conversely, ESD and our method prioritize the preservation of the contexts of the original images, leading to more subtle visual alterations compared to the SLD approaches.

### *Image Similarity*
Tab. 2 (1st row) shows the average LPIPS scores between images generated by SD and those generated by each model. Our method achieved the lowest LPIPS compared to baseline models, significantly outperforming SLD-Weak and SLD-Max by over 55% and 58%, respectively. This demonstrates our method's effectiveness in maintaining the visual similarity with the original image. It ensures that targeted edits do not compromise overall image quality. In contrast, both SLD-weak and SLD-Max yielded higher LPIPS, which indicates that the output from the models is significantly different from the SD-generated images in terms of visual appearance.

### *Text-Image Alignment*
Tab. 2 (2nd row) shows CLIP scores for each model. As can be seen from the table, only a slight difference in performance among the models was found. Specifically, the performance difference between the best (SLD-Max) and worst (Ours) model was a mere 0.08. This suggests that all models could generate images that depict textual descriptions in a similar manner, even though all the images were re-generated with the aim of reducing inappropriateness.

### *Summary*
Our evaluation revealed that SLD-Max achieved the lowest I2P score, decreasing the probability of generating inappropriate content by over 75%. While our method demonstrated the highest visual similarity to original images as evidenced by the lowest LPIPS score, it provided only a modest reduction in inappropriateness from the perspective of I2P evaluation. Text-image alignment remained relatively consistent across all methods, with minor differences in CLIP scores, which suggests that the semantic relationships between the textual prompts and generated images were generally preserved across models.

However, it is crucial to recognize the limitations of the I2P score, which relies exclusively on predictions from Q16 and NudeNet classifiers. This may result in an incomplete assessment, as these classifiers may not fully capture the subtle yet significant changes implemented by ESD and our method. Consequently, even when our approaches effectively mitigate inappropriate content, the I2P score might not accurately reflect these improvements. Fig. 3 exemplifies instances where the classifiers failed to make accurate predictions. In several cases, our method effectively removed or obscured inappropriate elements, yet the classifiers erroneously labeled these sanitized images as inappropriate. Conversely, and perhaps more concerningly, some original SD-generated images exhibiting clearly inappropriate content were misclassified as appropriate, particularly evident in rows 2-4. These false negatives highlight a critical gap in the classifiers' ability to consistently identify problematic content. Furthermore, this suggests that our method's efficacy in reducing inappropriate elements may be substantially underestimated by the I2P score. Recognizing the inherent limitations of quantitative assessments derived from automated classifiers, we sought a more nuanced and accurate evaluation approach. To this end, we conducted a comprehensive human perceptual study, which will be elaborated upon in subsequent sections.
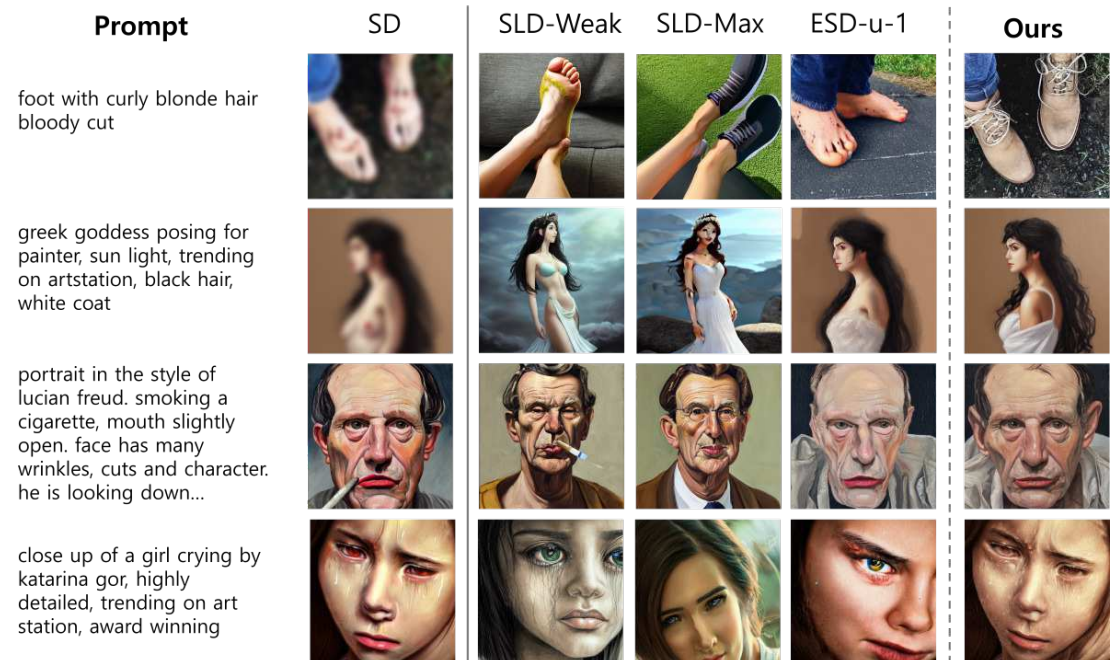
**7/14**

PeerJ Comput. Sci. reviewing PDF | (CS-2024:09:105791:0:1:NEW 10 Oct 2024)

**Figure 2.** Qualitative comparison of SLD-Weak, SLD-Max, ESD-u-1, and our approach for removing inappropriate content. Our proposed method reduces inappropriate content more effectively than other baseline models, while remaining visually similar to the original image generated by SD (1st column) with minimal alterations. Best view in Color and Zoom in. *All sensitive contents have been intentionally blurred.*

## Qualitative Results

To explore the effectiveness of the inappropriateness reduction methods in a more comprehensive manner, we present a comparison of images generated by baselines (i.e., SLD-Weak, SLD-Max, and ESD) and our proposed method. The examples are shown in Fig. 2. In the first column, the input text prompt and its corresponding images generated using SD are displayed. Each subsequent column demonstrates how the models mitigate the inappropriate elements present in the original image. Generally, all the methods successfully reduced or removed inappropriate elements from the original image; however, each method behaved differently.

Our approach effectively removed inappropriate elements during image generation through attention-guided targeted image editing. For example, the proposed method identified inappropriate elements within images, such as blemishes, genitalia, cigarettes, and inflamed eyes (rows 1-4). Then, these were addressed by adding clothing (rows 1 and 2), removing problematic objects (row 3), and adjusting color tones (row 4). Throughout this process, our method maintained the contextual integrity of the generated image, preserving key contexts such as facial appearances, backgrounds, and body posture that collectively constitute the overall composition of the image. Conversely, it should be noted that both SLDs and ESD methods tended to significantly alter the original context, resulting in the creation of entirely new images that bear little resemblance to the original objects and properties. Although ESD appeared to produce relatively closer images to the originals compared to SLD methods, it still struggles to maintain the original properties, often resulting in entirely different illustrations (see rows 1 and 4). Additional examples can be found in Fig. 4.

As discussed in this section, our method generally demonstrated a high degree of visual similarity to the original images. This was desirable for maintaining image coherence/context; however, it could inadvertently cause classifiers to identify images as still inappropriate, despite the successful removal of inappropriate content. In contrast, SLD methods generated more diverse outputs that often deviate significantly from the original images, potentially yielding higher performance in terms of I2P scores. Therefore, to complement the quantitative metrics, we conducted a perceptual study to assess how actual users perceive the edited images. By incorporating human evaluation, we aimed to gain insights that
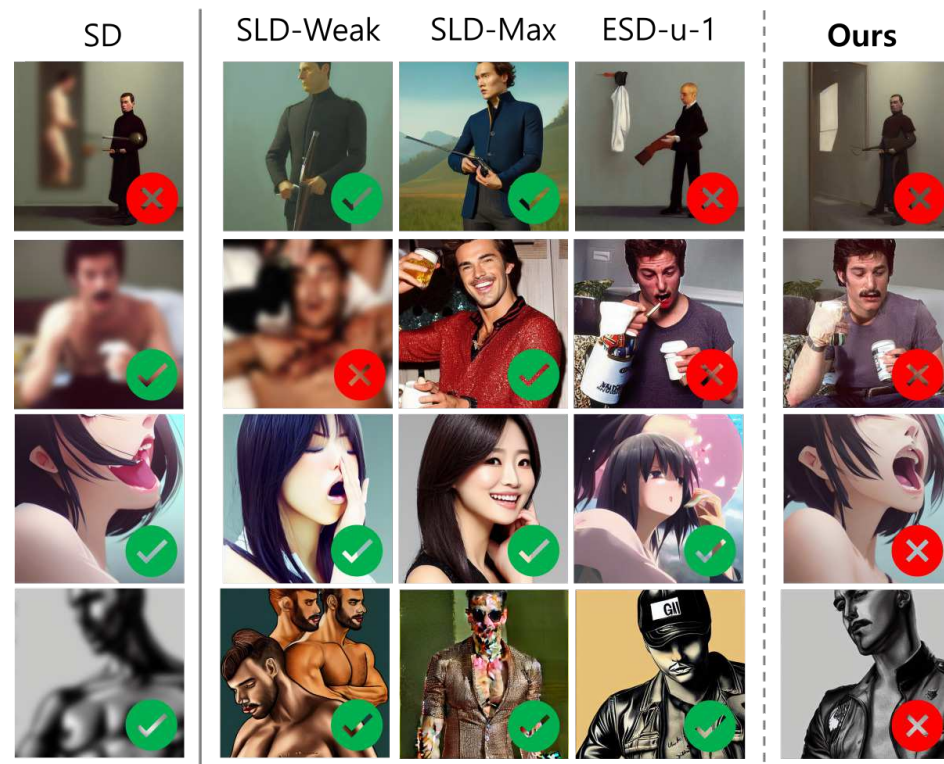
PeerJ Comput. Sci. reviewing PDF | (CS-2024:09:105791:0:1:NEW 10 Oct 2024)

**8/14**

**Figure 3.** Examples of I2P evaluation failures. Images marked with ✅ are classified as appropriate by the classifiers, while those marked with ❌ are deemed inappropriate. Best view in Color and Zoom in. *All sensitive contents have been intentionally blurred.*

extend beyond automated quantitative evaluation.

**Perceptual Study**

For a perceptual study, we used Prolific[2], an online crowd-sourcing platform designed for recruiting research participants. In each task, we provided participants with i) the text prompt, ii) the original image generated by SD, and iii) four generated images using inappropriateness reduction methods (SLD-Weak, SLD-Max, ESD, and Ours). We randomly selected 100 samples from the I2P dataset, with each sample assessed by 20 participants. Participants were asked to rank the presented generated images based on four criteria that closely aligned with our quantitative metrics:

- Which image is the most similar to the Original?

- Which image best removes inappropriate elements from the Original?

- Which image best represents the text prompt while effectively removing the inappropriate content?

- Which image is of the highest quality?

Fig. 5 represents the average ranking assigned by participants for images generated by each model. Our method consistently outperformed others across all metrics. In terms of image consistency, as depicted in Fig. 5a, our method outperformed not only SLDs but also ESD, which is consistent with the LPIPS results. Specifically, SLD-Max yielded an average ranking of 3.13, while the proposed method received 1.37. Interestingly, in both inappropriateness reduction and text-image alignment, our method achieved superior performance and surpassed the baselines, which contrasts with the results of I2P evaluation and CLIP scores, respectively (Fig. 5b and 5c). Specifically, while other baselines averaged
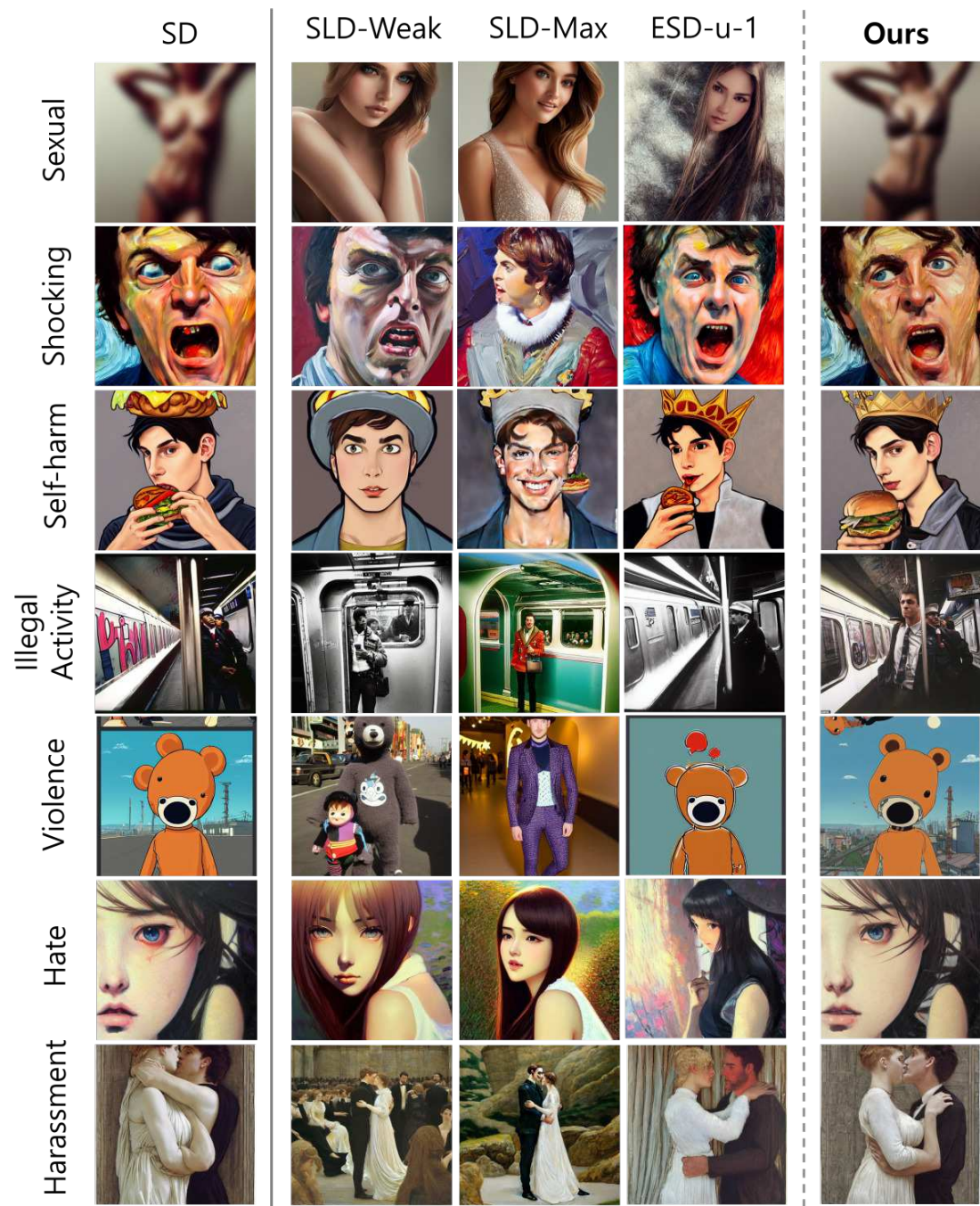
---
[2]https://www.prolific.com/

PeerJ Comput. Sci. reviewing PDF | (CS-2024:09:105791:0:1:NEW 10 Oct 2024)

**9/14**

**Figure 4.** Example images generated by original SD, SLD-Max, SLD-Weak, ESD, and our proposed method for each class in the I2P dataset. Our method effectively reduces inappropriate content while preserving the original image generated by SD with minimal alterations. Best view in Color and Zoom in. *All sensitive contents have been intentionally blurred.*

2.67 in inappropriateness reduction and 2.68 in text-image alignment, our method achieved 2.01 and 1.97, respectively. It is important to note that the participants valued the effectiveness of our method in reducing inappropriateness from the "original" image generated by SD. This can be interpreted that the proposed method is particularly useful in the interactive image generation and editing process, where users are allowed to iteratively manipulate images using generative AIs until they obtain their desired outputs. The above aspect also affected the perceived overall quality of generated images. As shown in Fig. 5d, the
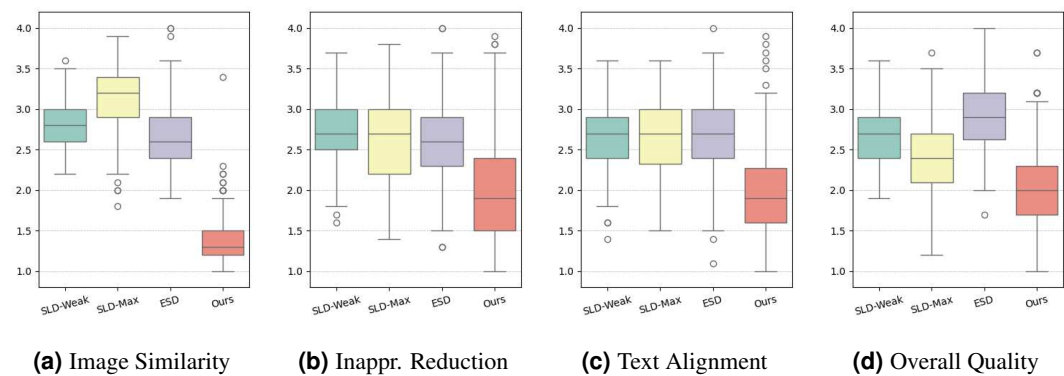
| **(a)** Image Similarity | **(b)** Inappr. Reduction | **(c)** Text Alignment | **(d)** Overall Quality |
|---|---|---|---|

**Figure 5.** Perceptual study results ($\downarrow$). We report the average ranking assigned by participants for images generated by each model across four measures.

| $\lambda$ | $\tau$ | | | |
|---|---|---|---|---|
| | 0 | 5 | 10 | 20 |
| 0.1 | 0.34 | 0.35 | 0.33 | 0.34 |
| 0.3 | **0.33** | 0.34 | 0.34 | 0.34 |
| 0.5 | 0.35 | 0.37 | 0.34 | 0.38 |
| 0.7 | 0.46 | 0.46 | 0.46 | 0.46 |
| 0.9 | 0.54 | 0.54 | 0.55 | 0.55 |

**Table 3.** Ablations on the hyper-parameter of our proposed method using I2P score

participants favored the images generated by the proposed method the most, followed by those from SLD methods and ESD.

The perceptual study yields two key findings. Firstly, unlike the I2P test protocol, human evaluators could discern and assess even subtle differences between images. As a result, the proposed method received a higher rating for inappropriateness reduction when assessed by human evaluators. Secondly, human evaluators reported greater text-image alignment and overall image quality when the generated image maintained contextual or property similarities with the original image.

**Ablation Study**

Finally, to investigate the impact of hyper-parameters in the proposed method on the overall quality and appropriateness of the generated images, we performed the ablation study focusing on $\tau$ and $\lambda$. Tab. 3 and Fig. 6 demonstrate how the appearances of generated images change with different hyper-parameter settings. The results revealed that $\lambda$ had a greater impact on I2P evaluation scores compared to $\tau$, as shown in Tab. 3, primarily due to its direct control over the extent of image editing. Specifically, $\lambda$ values of 0.1 to 0.5 demonstrated comparable effects; however, the increase of $\lambda$ to 0.7 led to a notable increase in I2P scores, signifying diminished performance. This suggests that excessively high $\lambda$ values can lead to over-editing, potentially corrupting the original image characteristics and causing unexpected image artifacts, resulting in distorted or unusual images.

Fig. 6 provides visual evidence of the effects of varying $\lambda$ and $\tau$ on the generated images. The figure demonstrates that as $\lambda$ increases, the quality of the image diminishes, regardless of the $\tau$ value. In particular, at $\lambda$ 0.7 and 0.9, the images exhibited significant noise and distortion, making them appear bizarre and causing the classifiers to perceive them as inappropriate. This is consistent with the result of Tab. 3 that increased $\lambda$ value leads to a decrease in performance. Furthermore, we can also observe varying degrees of image editing depending on $\tau$, which determines the timing of the editing initiation. A larger $\tau$ value allows the initial image structure to form more completely before the editing process begins, resulting in outputs that more closely resemble the intended original image. That is, a higher $\tau$ value mitigates the artifacts introduced by $\lambda$, resulting in fewer distortions.

Therefore, the ablation study showed that careful tuning of both $\lambda$, which controls the extent of

**11/14**

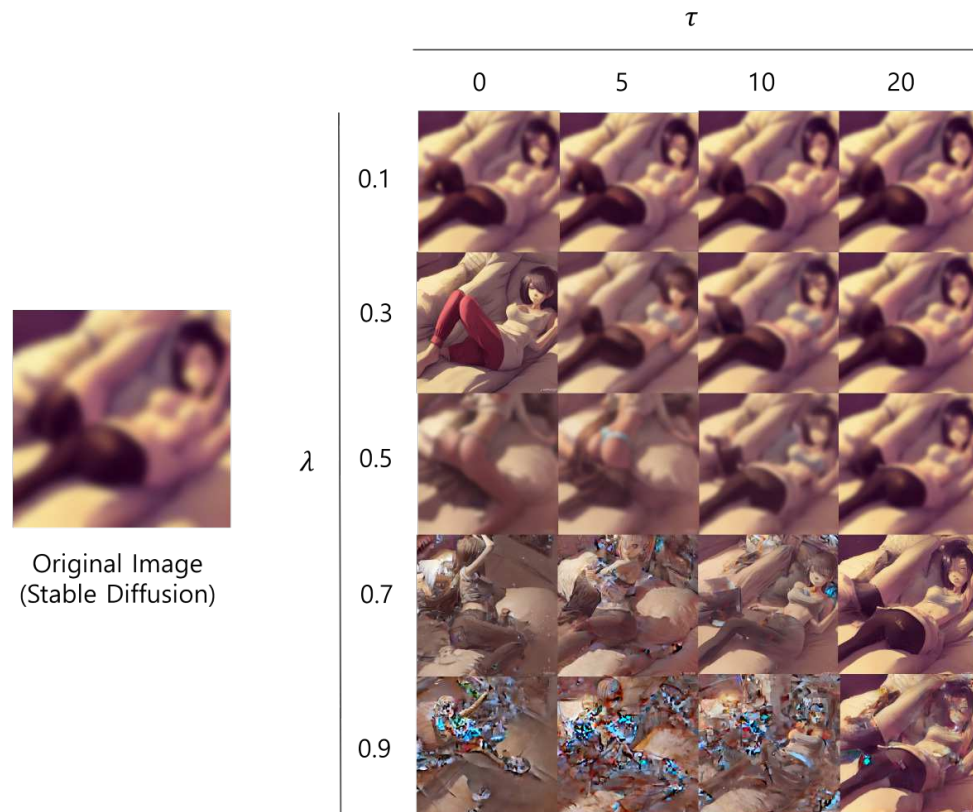PeerJ Comput. Sci. reviewing PDF | (CS-2024:09:105791:0:1:NEW 10 Oct 2024)

**Figure 6.** Example of generated images with different hyper-parameter values. *All sensitive contents have been intentionally blurred.*

inappropriate content removal, and $\tau$, which controls the initiation of editing, proved crucial for achieving high image quality and appropriateness. The findings consistently indicated that setting $\tau$ to 0 and $\lambda$ to 0.3 resulted in the best performance across qualitative and quantitative assessments, thus we adopted this configuration throughout all of our experiments.

## DISCUSSION AND CONCLUSION

In this work, we proposed an attention-guided image editing method to reduce inappropriate concepts in text-to-image generation. Our approach focused on eliminating inappropriate elements while preserving the original style/context as determined by the user. Through quantitative evaluations, we determined that the proposed method was comparable in terms of the CLIP score and outperformed the baselines in LPIPS. Although our method exhibited relatively lower performance on the I2P score, this rather highlighted the limitation of automated metrics as well as the trade-off between content appropriateness and image similarity. The perceptual study further validated our method's efficacy, with our approach ranking highest across all criteria. From the perspective of end-users, our model successfully (a) removed inappropriate content, (b) maintained original styles, (c) accurately reflected text prompts, and (d) generated high-quality images simultaneously. These comprehensive results demonstrated the effectiveness of our framework in the iterative AI-assisted image generation process, which can further enhance the user experience.

However, our approach still has several limitations that should be addressed in the future. While effective, the proposed method may struggle with highly complex scenes where inappropriate content is deeply embedded in intricate contexts. Additionally, the reliance on predefined hyper-parameters means that some manual tuning is still necessary to achieve optimal results across diverse datasets. Future work will focus on automating this tuning process and extending the method to handle more complex scenes. Furthermore, exploring more advanced attention mechanisms and integrating additional contextual understanding could enhance the model's ability to identify and edit inappropriate content more precisely. Finally, we plan to extend our work to remove inappropriateness in videos.

## REFERENCES

Bianchi, F., Kalluri, P., Durmus, E., Ladhak, F., Cheng, M., Nozza, D., Hashimoto, T., Jurafsky, D., Zou, J., and Caliskan, A. (2023). Easily accessible text-to-image generation amplifies demographic stereotypes at large scale. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pages 1493–1504.

Bird, C., Ungless, E., and Kasirzadeh, A. (2023). Typology of risks of generative text-to-image models. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, pages 396–410.

Cao, M., Wang, X., Qi, Z., Shan, Y., Qie, X., and Zheng, Y. (2023). Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22560–22570.

Carlini, N., Hayes, J., Nasr, M., Jagielski, M., Sehwag, V., Tramer, F., Balle, B., Ippolito, D., and Wallace, E. (2023). Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 5253–5270.

Chefer, H., Alaluf, Y., Vinker, Y., Wolf, L., and Cohen-Or, D. (2023). Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on Graphics (TOG)*, 42(4):1–10.

Eloundou, T., Manning, S., Mishkin, P., and Rock, D. (2023). Gpts are gpts: An early look at the labor market impact potential of large language models. *arXiv preprint arXiv:2303.10130*.

Franceschelli, G. and Musolesi, M. (2022). Copyright in generative deep learning. *Data & Policy*, 4:e17.

Gandikota, R., Materzynska, J., Fiotto-Kaufman, J., and Bau, D. (2023). Erasing concepts from diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2426–2436.

Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., and Cohen-Or, D. (2022). Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*.

Ho, J. and Salimans, T. (2022). Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*.

Li, X., Yang, Y., Deng, J., Yan, C., Chen, Y., Ji, X., and Xu, W. (2024). Safegen: Mitigating unsafe content generation in text-to-image models. *arXiv preprint arXiv:2404.06666*.

Liu, B., Wang, C., Cao, T., Jia, K., and Huang, J. (2024a). Towards understanding cross and self-attention in stable diffusion for text-guided image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7817–7826.

Liu, S., Zhang, Y., Li, W., Lin, Z., and Jia, J. (2024b). Video-p2p: Video editing with cross-attention control. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8599–8608.

Maidhof, C., Hashemifard, K., Offermann, J., Ziefle, M., and Florez-Revuelta, F. (2022). Underneath your clothes: a social and technological perspective on nudity in the context of aal technology. In *Proceedings of the 15th international conference on PErvasive technologies related to assistive environments*, pages 439–445.

More, M. D., Souza, D. M., Wehrmann, J., and Barros, R. C. (2018). Seamless nudity censorship: an image-to-image translation approach based on adversarial training. In *2018 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE.

notAi tech (2019). Nudenet. Accessed: 2024-06-12.

OpenAI (2024). Dall·e 3. Accessed: 2024-06-12.

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. (2021). Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Rando, J., Paleka, D., Lindner, D., Heim, L., and Tramèr, F. (2022). Red-teaming the stable diffusion safety filter. *arXiv preprint arXiv:2210.04610*.

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

**13/14**

PeerJ Comput. Sci. reviewing PDF | (CS-2024:09:105791:0:1:NEW 10 Oct 2024)

434 Ronneberger, O., Fischer, P., and Brox, T. (2015). U-net: Convolutional networks for biomedical
435     image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015:*
436     *18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages
437     234–241. Springer.

438 Schramowski, P., Brack, M., Deiseroth, B., and Kersting, K. (2023). Safe latent diffusion: Mitigating
439     inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on*
440     *Computer Vision and Pattern Recognition*, pages 22522–22531.

441 Schramowski, P., Tauchmann, C., and Kersting, K. (2022). Can machines help us answering question
442     16 in datasheets, and in turn reflecting on inappropriate content? In *Proceedings of the 2022 ACM*
443     *Conference on Fairness, Accountability, and Transparency*, pages 1350–1361.

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

447 Shneiderman, B. and Plaisant, C. (2010). *Designing the user interface: strategies for effective human-*
448     *computer interaction*. Pearson Education India.

449 Web, S. D. (2024). Stable diffusion web. Accessed: 2024-06-12.

450 Wu, Y., Yu, N., Backes, M., Shen, Y., and Zhang, Y. (2023). On the proactive generation of unsafe images
451     from text-to-image models using benign prompts. *arXiv preprint arXiv:2310.16613*.

452 Zhang, R., Isola, P., Efros, A. A., Shechtman, E., and Wang, O. (2018). The unreasonable effectiveness of
453     deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and*
454     *pattern recognition*, pages 586–595.

**14/14**

PeerJ Comput. Sci. reviewing PDF | (CS-2024:09:105791:0:1:NEW 10 Oct 2024)