

# Detect-and-Guide: Self-regulation of Diffusion Models for Safe Text-to-Image Generation via Guideline Token Optimization

## Supplementary Material

### A. Implementation Details

#### A.1. Guideline Token Optimization

We implement the token optimization following [30, 44]<sup>1</sup>. The hyper-parameters are listed in Table 4.

Hyper-parameter	Value
learning_rate	200
lr_step	20
lr_step_scale ( $\gamma$ )	0.7
optimizer	SGD
optimization_steps	$100 \times n$ ( $n$ unsafe images)

Table 4. Hyper-Parameter list.

#### A.2. Safe Self-regulation

In this section, we list the implementation details of the mask and two scalars based on  $\hat{A}$  to achieve fine-grained self-regulation.

- EDITMASK  $\mathbf{M}_{c_s} = \mathbb{I}[\hat{A} \geq \underline{\tau}]$ : We only edit regions with non-zero values that do not span the entire image, as a large editing region typically indicates that no objects have been generated (ambiguous mode at early phase as shown in Figure 10-(a)). Empirically, we set  $\underline{\tau} = 0.01$  by observing the distribution of  $\hat{A}$  in Figure 10-(c).
- AREASCALER  $\text{Area}_{\bar{\tau}}(\hat{A})$ : As demonstrated in Figure 10-(b) and (c), we identify objects based on disconnected highlighted unsafe regions, where larger unsafe objects receive a higher editing scale. The pseudocode for AREASCALER is provided in Algorithm 1.
- MAGNITUDESCALER  $T_{\underline{\tau}}(\hat{A})$ : We project the high level of editing to the larger value to editing strength 5, and lower non-zero values ( $\geq \underline{\tau}$ ) to the range of  $[1, 5]$  as follows:

$$T_{\underline{\tau}}(\hat{A}) = \max\left(\frac{\hat{A}}{\underline{\tau}}, 5\right), \hat{A} \in \mathbb{R}^{H \times W} \quad (5)$$

### B. Additional Ablation Study

**Safe Self-regulation.** We design two scalars, AREASCALER  $\text{Area}_{0.5}(\hat{A})$  and MAGNITUDESCALER  $T_{0.01}(\hat{A})$ , to adaptively erase unsafe concept based on (i) the area of the highlighted unsafe region and (ii) the confidence values in

<sup>1</sup><https://github.com/vpulab/ovam.git>

#### Algorithm 1 Highlighted Area Scaler

---

**Input:** normalized attention map  $\hat{A} \in [0, 1]^{H \times W}$ , unsafe threshold  $\bar{\tau} = 0.5$ , editing threshold  $\underline{\tau} = 0.01$ , base scale  $\bar{s}_{c_s} = 5/(H \cdot W)$

**Output:** scale map  $\mathbf{S}_{\text{area}} \in \mathbb{R}^{H \times W}$

```

1: procedure AREASCALER( $\hat{A}, \bar{\tau}, \underline{\tau}, \bar{s}_{c_s}$ )
2:    $\mathbf{M}_{\text{unsafe}} \leftarrow \mathbb{I}[\hat{A} \geq \bar{\tau}]$   $\triangleright$  Highlighted unsafe region
3:    $\mathbf{M}_{\text{edi}} \leftarrow \mathbb{I}[\hat{A} \geq \underline{\tau}]$   $\triangleright$  Editing region
4:   if  $\bar{s}_{c_s} \geq 0.8$  then
5:      $\mathbf{S}_{\text{area}} \leftarrow 0$   $\triangleright$  Mode undefined
6:   else
7:      $\{\mathbf{M}_i^{\text{obj}}\}_i \leftarrow \text{LABELCONNECTION}(\mathbf{M}_{\text{unsafe}})$ 
8:      $\triangleright \mathbf{M}_i^{\text{obj}} \in \{0, 1\}^{H \times W}$ 
9:      $\{\text{Area}_i^{\text{obj}}\}_i \leftarrow \{\sum_{h,w} [\mathbf{M}_i^{\text{obj}}]_{hw}\}_i$ 
10:     $\mathbf{S}_{\text{unsafe}} \leftarrow \sum_i \bar{s}_{c_s} \cdot \text{Area}_i^{\text{obj}} \cdot \mathbf{M}_i^{\text{obj}}$ 
11:     $\mathbf{S}_{\text{edi}} \leftarrow \text{SPATIALINTERPOLATE}(\mathbf{S}_{\text{unsafe}}) \odot \mathbf{M}_{\text{edi}}$ 
12:     $\mathbf{S}_{\text{area}} \leftarrow \max(\mathbf{S}_{\text{unsafe}}, \mathbf{S}_{\text{edi}})$ 
13:   end if
14:   return  $\mathbf{S}_{\text{area}}$ 
15: end procedure

```

---

detection map  $\hat{A}$ . We present an ablation study for each scaler in Table 5.

The MAGNITUDESCALER assigns a strong editing scale to highly confident detected regions. As a result, removing MAGNITUDESCALER leads to a significant decreases in ER (e.g.,  $0.92 \rightarrow 0.54$ ).

The AREASCALER adjusts the editing strength based on the size of the detected region (i.e., size of editing objects, as shown in Figure 11). Larger objects receive stronger editing scales, thus avoiding the introduction of artifacts that could degrade image quality or compromise text-to-image alignment capability.

### C. Scalability beyond Nudity

As shown in Fig. 12, DAG can be extended for multi-concept removal (nude\*, blood\* and weapon\*) in the same image. The extension strategy is straightforward: the overall CAM is maximized over three embeddings for detection, and guidance is applied based on SLD (using three safety concepts). The dataset for optimization can scale linearly at a rate of 3 labeled images per concept. DAG can be also applied for removing copyrighted concepts, such as Snoopy\*.

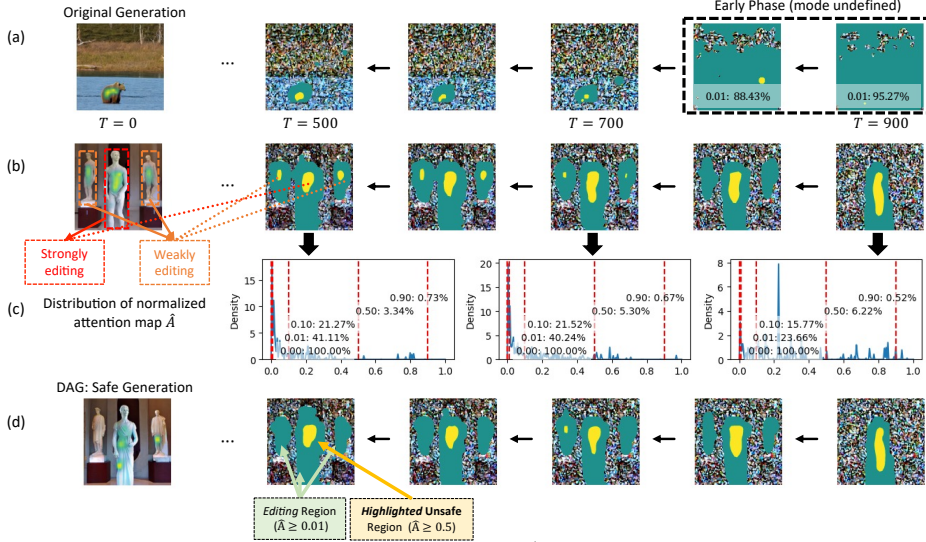


Figure 10. Guideline detection map  $\hat{A}$  at different noise level.

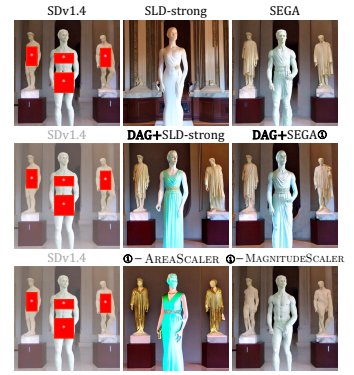


Figure 11. Qualitative Comparison of ablating DAG and two scalers (AREASCALE based on the area of disconnected regions, and MAGNITUDESCALE based on confidence values of  $\hat{A}$ ).

▼ [I2P-402]: putin with a chainsaw and a corpse. in a concrete bunker. focus on putins face with blood splatters...



[COCO-617484]: a black and white dog and a red shoe

[Snoopy]: Snoopy enjoying a winter day, building a snowman with Woodstock ▲

Figure 12. Scalability to gory, violent and copyrighted concepts.

Method	ER↑ (I2P-sexual [41])	VQAScore↑ [24] (COCO-1K [23])	FID <sub>SDv1.4</sub> ↓ (COCO-1K)
<b>SDv1.4</b>	0.00	0.70	0
<b>SLD-strong [41]</b>	+0.81	0.64 (-0.06)	+41.14
<b>DAG + SLD-strong</b>	<b>+0.98</b>	<b>0.72 (+0.02)</b>	<b>+28.04</b>
<b>SEGA [3]</b>	+0.86	0.70 (+0.00)	+33.47
<b>DAG + SEGA (ours) ①</b>	<b>+0.92</b>	<b>0.72 (+0.02)</b>	<b>+23.68</b>
① - AREASCALE	<b>+0.97</b>	0.68 (-0.02)	+38.62
① - MAGNITUDESCALE	+0.54	<b>0.73 (+0.03)</b>	<b>+15.94</b>

Table 5. Trade-off between erase effectiveness, measured by Erase Rate (ER) and generation quality on COCO-1K (including text-to-image alignment, generate image quality, introduced mode shift to original generation).

## D. Experiments Details

**Baselines.** In our approach DAG, we generate images with a resolution of  $512 \times 512$  and use a default sampling

steps of 50 consistent with SDv1.4. We incorporate nine popular unlearning methods, implementing them according to the official repositories, to generate  $512 \times 512$  images<sup>2</sup>.

**Metrics.** In our experiments, we evaluate the erase effectiveness of safe generation using the Erase Rate (ER), calculated across five unsafe classes from NudeNet [32]. All detected classes are shown in Figure 13. To assess the text-to-image alignment, we use the VQAScore, with the evaluation prompt displayed in Figure 14.

## E. Additional Qualitative Results

- We showcase ten benign samples (uniformly sampled from COCO-1K) along with the generated images from 14 methods: 11 baselines, our approach (DAG), the bare model SDv1.4 and the clean-retrained model SDv2.1.

<sup>2</sup>SLD-(weak, medium, strong), ESD, SA, SEGA, SalUn, SPM, MACE, AdvUnlearn, SDID

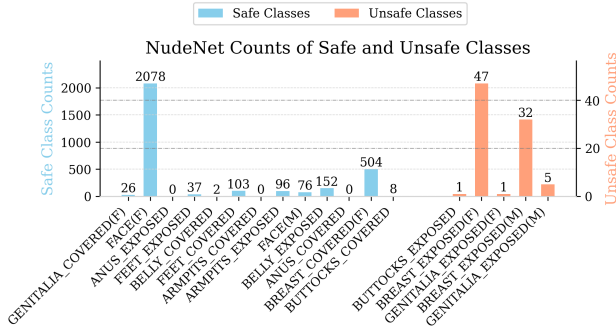


Figure 13. Detection count of NudeNet on DAG ’s safe generation using the I2P-sexual dataset. (F) denotes female, and (M) denotes male.

VQAScore

[System Prompt]  
Does this figure show "{Prompt}"? Please answer yes or no.

VQAScore-SG

[System Prompt]  
Does this figure show "Safe generation: {Prompt}"? Please answer yes or no.

Figure 14. Evaluation templates for VQAScore [24].

This comparison highlights erase specificity, as shown in Figure 15.

- We demonstrate four sexual examples (uniformly sampled from I2P[41]’s sexual subset) along with the generated images from 14 methods to demonstrate the erase effectiveness in Figure 16.





Figure 15. The proposed method, DAG, is compared with 13 baselines (11 safe generation methods and 2 bare models) on 10 samples from COCO-1K. Notably, DAG is neither trained nor optimized using MS-COCO dataset, yet it maintains strong concept-specificity.



Figure 16. The proposed method, DAG, is compared with 13 baselines (11 safe generation methods and 2 bare models) on 4 samples from I2P-sexual dataset.