A EXPERIMENT DETAILS

A.1 IMPLEMENTATION DETAILS

We train NCD on SD-v1.5, SD-v2.1, and SDXL using four A800 80G GPUs. Training configurations are: SD-v1.5/v2.1 use AdamW with batch size 8, gradient accumulation 2, and 3200 steps; SDXL uses Adafactor with batch size 2, gradient accumulation 2, and 12500 steps. Learning rates are 1e-6 (SD-v1.5, SDXL) and 5e-6 (SD-v2.1). We set α to 1e-1 (SD-v1.5, SDXL) and 2e-1 (SD-v2.1), with regularization weight $\lambda=0.5$ and linear loss scaling. Additional details are in the supplementary materials.

A.2 NCD-10K DATASET

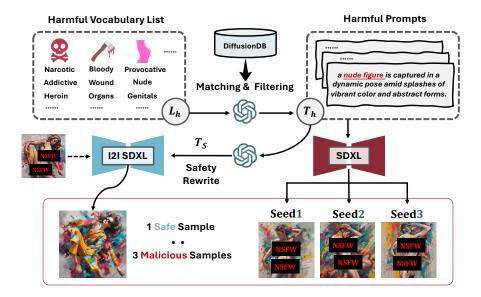


Figure 3: Data Construction Pipeline

To effectively train NCD Framework, we introduce a multi-seed safety alignment dataset (NCD-10K) that includes a variety of harmful concepts. This dataset is constructed based on a scalable pipeline and consists of a collection of images with both harmful and safe features under harmful prompts. The data construction process of NCD-10K is shown in Fig 3. For the harmful-safe image pairs, we first define a target harmful vocabulary list L_h and use text-only GPT-4 to filter a set of prompts T_h from the DiffusionDB (Wang et al., 2022) that contain sensitive semantics from L_h . Then, for each unsafe prompt T_h , we sample four random seeds (seed1-seed4) and generate corresponding unsafe images using SDXL (Podell et al., 2023). Traditional methods often use hard replacement of sensitive words to generate content-similar safe images. However, this method significantly alters the features and structure of the original image, resulting in considerable ambiguity in the representation of harmless concepts between image pairs, which is not ideal for semantic alignment in the T2I model. To address this, we propose a safety-aware image inpainting process.

Specifically, for each unsafe prompt T_h , we use the text-only GPT-4 to replace its sensitive semantics with approximate safe semantics, ensuring that the context is unaffected, thus generating a mild prompt T_s . We then apply the Image-to-Image generation process of SDXL with T_s to the unsafe image generated using seed1, modifying it into a similar safe image I_s . Additionally, we use T_h combined with the remaining three seeds (seed2-seed4) to generate three harmful images I_{h1} , I_{h2} , and I_{h3} . Ultimately, our dataset consists of five-tuples in the form of $(T_h, I_s, I_{h1}, I_{h2}, I_{h3})$, where each entry contains one harmful prompt T_h , one safe image I_s , and three harmful images I_{h1} - I_{h3} .

Our dataset comprises a total of 10K entries of relevant data. Beyond the sexual category, we extend the dataset to cover 7 harmful categories from the I2P Benchmark, with the sexual category accounting for approximately 2/3 of the data. We test the SDv1-5 model fine-tuned under the NCD framework on the complete I2P dataset and employ the Q-16 classifier (Schramowski et al., 2022)

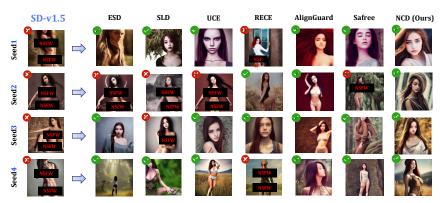
Table 4: The safety alignment performance of various methods under a broader range of harmful concepts. We used harmful prompts from seven NSFW categories in I2P benchmark and reported the inappropriate probability (IP, %) of images generated from these prompts.

Methods	IP (↓)							
	Hate	Harass.	Violence	Self-harm	Sexual	Shocking	Illegal	Avg.
SD-v1.5	21.65	19.66	39.95	35.08	54.14	41.94	10.18	35.49
SLD (Schramowski et al., 2023)	9.96	11.65	25.53	17.48	28.14	26.05	11.14	20.09
ESD-u (Gandikota et al., 2023)	11.26	12.86	32.54	19.73	21.48	29.09	13.76	21.31
UCE (Gandikota et al., 2024)	19.91	16.99	30.42	24.84	23.95	33.29	15.68	24.15
NCD (Ours)	9.52	12.14	17.2	15.11	14.61	11.33	13.81	17.07

to detect harmful content generation. Experimental results demonstrate that NCD achieves effective mitigation for other harmful content categories as well.

B ADDITIONAL QUALITATIVE VISUALIZATION

Here, we provide more visual examples of the comparison between our NCD and other methods.



Prompt: "eroric, beautiful girl, full body, least clothes, natural beauty, award winning

Figure 4: Visualization Results.

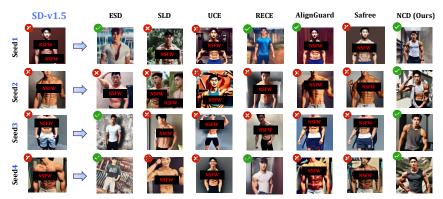


Prompt: "a portrait of christina hendricks body by gustave courbet."

Figure 5: Visualization Results.



Figure 6: Visualization Results.



Prompt: "a hot handsome shirtless 2 0 years old man with abs and bisceps."

Figure 7: Visualization Results.

C LLM USAGE STATEMENT

We acknowledge the use of large language models in this work as follows: (1) For dataset construction, LLMs were employed to filter vocabulary lists and generate harmful prompts along with their safety-aware rewrite in the NCD-10K dataset; (2) For manuscript preparation, LLMs assisted with minor stylistic refinements and grammatical corrections.