

Focus-N-Fix: Region-Aware Fine-Tuning for Text-to-Image Generation

Supplementary Material

This appendix is structured as follows:

- In Appendix A, we provide additional related work.
- In Appendix B, we provide experimental details.
- In Appendix C, we study the performance of our proposed model using objective metrics.
- In Appendix D, we conduct ablation study for hyper-parameters used in our region-aware fine-tuning.
- In Appendix E, we show sample heatmaps for over-sexualization (safety) reward fine-tuning from the training and test datasets.
- In Appendix F, we provide additional results for mitigating over-sexualization, artifacts, and violence. We also provide qualitative results on forgetting.
- In Appendix G, we show that our method can be applied to other diffusion models besides Stable Diffusion v1.4.
- In Appendix H, we provide a performance comparison between our proposed method, *Focus-N-Fix* and a popular concept editing method, UCE [17].

A. Additional Related Work

Evaluation and Rewards for Image Generation. Early works proposed automated metrics for image evaluation, like Fréchet Inception Distance (FID) [22], Inception Score (IS) [44], and Learned Perceptual Image Patch Similarity (LPIPS) [60]. To evaluate vision-language alignment, CLIPScore [21] has been commonly used to measure the similarity of the image and prompt. However, these metrics still fall short in reflecting human preferences. More recent work has introduced higher quality datasets such as HPSv2 [53], PickScore [27] and ImageReward [55] that collect human preference annotations to guide image evaluation. RichHF [32] further enriches the feedback signal related to unsatisfactory image regions and prompt tokens missing from images. Additionally, with the rapid development of large vision language models (LVLM), some current works leverage LVLMs to simulate human rewards. Among them, DreamSync [49], TIFA [25], VIEScore [28], LLMscore [36] utilize Visual Question Answering (VQA) tasks to quantitatively assess image generation qualities. However, these evaluation methods are highly dependent on the performance of LVLMs.

B. Experimental Details

In this section, we provide detailed experiment settings of our proposed method and the baselines. We primarily implement the methods using Stable Diffusion (SD) v1.4. We use a sampling process with 50 steps and a classifier-free guidance weight of 7.5. The resolution of the generated

images is 512×512 pixels.

Reward Fine-tuning Settings. We fine-tune SDv1.4 using LoRA parameters with a rank of 64 and optimize the parameters using the AdamW optimizer [34]. We adopt an initial learning rate of $\eta_0 = 2 \times 10^{-5}$ and a square root decay schedule, where the learning rate at training set i is $\eta_i = \eta_0/\sqrt{i}$. The scale of the regional constraint loss in Equation 2 is $\beta = 0.001$ for artifact reward fine-tuning and $\beta = 5e - 4$ for over-sexualization (safety) reward fine-tuning. Utilizing 4 TPU (v5p) units, the fine-tuning was completed in 8-10 hours

Reward Guidance Settings. For the reward guidance baseline, we employed a reward guidance scale of $\lambda = 2.0$. To avoid overly large modifications to some image samples that lead to distortion, we apply L-2 norm gradient clipping with a threshold of 2.0. We add guidance starting from step 10 out of 50 sampling steps in total.

Safe Latent Diffusion Settings. We used a safety scale of 500 and a safety threshold value of 0.03.

C. Objective Metrics

The results presented in Table 1 and Table 5, obtained through human evaluation, demonstrate the superiority of our region-aware fine-tuning method over DraFT in improving the targeted quality aspect during fine-tuning with the corresponding reward model, while minimizing degradation in other aspects of image quality. To further strengthen our claim that region-aware fine-tuning only modifies problematic regions while preserving the rest of the image, we compute three image similarity metrics : Peak Signal-to-Noise Ratio (PSNR), SSIM [51] and LPIPS. These metrics are calculated by comparing the original SD v1.4 output with the results from either DraFT or our method, *Focus-N-Fix*, for both the full evaluation set and the smaller subset used in human evaluation. The results in Tables 3 and 4, across both safety and artifact fine-tuning experiments, show that the global image-level changes are notably smaller when fine-tuning with *Focus-N-Fix*.

D. Ablation Study

The heatmap constraint (β) in our proposed method *Focus-N-Fix* regulates the extent to which image generations from the fine-tuned model differ from those of the pre-trained

Safety Reward Fine-Tuning	Full Evaluation Set (419 Prompts)			Human Evaluation Set (100 Prompts)		
Method / Objective Metrics	PSNR (\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)	PSNR (\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
SD v1.4 vs DraFT	13.07	0.41	0.59	13.34	0.42	0.60
SD v1.4 vs Focus-N-Fix	22.30	0.80	0.18	23.64	0.83	0.16

Table 3. Metrics estimating image-Level Changes from the Original Stable Diffusion v1.4 Model in DraFT and Our Method, Focus-N-Fix, for Safety Fine-Tuning

Artifact Reward Fine-Tuning	Evaluation Set (HPDv2 Eval + Parti Prompts)			Human Evaluation Set (100 Prompts)		
Method / Objective Metrics	PSNR (\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)	PSNR (\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
SD v1.4 vs DraFT	15.39	0.55	0.45	15.04	0.54	0.45
SD v1.4 vs Focus-N-Fix	21.61	0.78	0.18	21.04	0.78	0.19

Table 4. Metrics estimating image-Level Changes from the Original Stable Diffusion v1.4 Model in DraFT and Our Method, Focus-N-Fix, for Artifact Fine-Tuning

model. When $\beta = 0$ our method simplifies to DraFT. In Figure 8, we illustrate the various improvements in the safety reward score as β is changed, while monitoring the perceptual changes in the image, measured by the objective metrics PSNR, SSIM, and LPIPS. Our aim is to enhance the safety score while ensuring that image generation from the fine-tuned model closely resembles the original image. To achieve this, we set thresholds for the objective metrics (PSNR > 20 , SSIM > 0.75 , and LPIPS < 0.2) and conduct extensive visual inspections, ultimately selecting $\beta = 5e-4$ as our chosen hyper-parameter. We repeated the analysis for artifact reward-finetuning experiment, ultimately selecting $\beta = 1e-3$ for the results in the paper.

E. Heatmaps in *Focus-N-Fix*

Focus-N-Fix leverages heatmaps during fine-tuning to detect and address problematic regions within the training data. This targeted strategy enables us to refine the model’s performance in these areas. Notably, during inference, despite the absence of heatmaps, the fine-tuned model shows a significant quality improvement. It achieves this by applying corrections in the problematic regions based on the adjustments learned during fine-tuning. Consequently, the proposed model consistently delivers image generations with a marked improvement in the problematic regions, demonstrating its ability to generalize effectively without relying on heatmaps for guidance. In Fig. 9, we present the safety reward heatmaps generated using a single prompt from both the training and evaluation datasets, along with the improvements achieved by *Focus-N-Fix* compared to the original images generated by the pre-trained SD v1.4.

F. Additional Results

In this section, we present additional qualitative results comparing various methods, including the pre-trained Stable Diffusion (SD) v1.4, Safe Latent Diffusion (SLD), Reward Guidance (with and without Region Constraint (RC)),

DRaFT, and our proposed method, *Focus-N-Fix*. We show examples of these methods mitigating over-sexualization, artifacts and violence contents. Additionally, we highlight sample failure cases of our *Focus-N-Fix* method and present a study on catastrophic forgetting, comparing our proposed fine-tuning approach, *Focus-N-Fix*, with DRaFT fine-tuning.

F.1. Over-sexualization

Fig. 10 shows results of the compared methods on addressing over-sexual contents in the generated images. It can be observed that SLD and DRaFT often generate images with different structure given the same prompts after fine-tuning to mitigate over-sexualization. Consequently, they may introduce other issues such as text-image misalignment and artifacts. For example, in Fig. 10 where the prompts are “A garden gnome wearing a bikini”, “A beautiful woman from behind walking on the beach”, and “Men at the beach” both SLD and DraFT outputs mitigate over-sexualization but at the expense of significantly altering the background regions and overall image structure compared to the original image generated by SD v1.4. In the example with prompt “portrait of a surreal reclining female body made of liquid forms, ink in water surreal paint swirls pastel colors” in Fig. 10 both SLD and DraFT reduce over-sexualization, but modify the original image in such a way that it reduces the text-to-image alignment with “female body” part of the prompt. DraFT and SLD methods can also introduce additional artifacts apart from changing the image structure as in example from Fig. 10 with prompt “Men at the beach”. Moreover, Reward Guidance may not have consistent performance across samples. In some cases, it may not effectively address the issues of concern (e.g., in the examples from Fig. 10 with prompts “A garden gnome wearing a bikini”, “A woman getting dressed in the closet, cinematic.”), while in other cases it may change the image structure (e.g., in the example from Fig. 10 with prompt “A curvy woman riding a horse.”). Even with the region con-

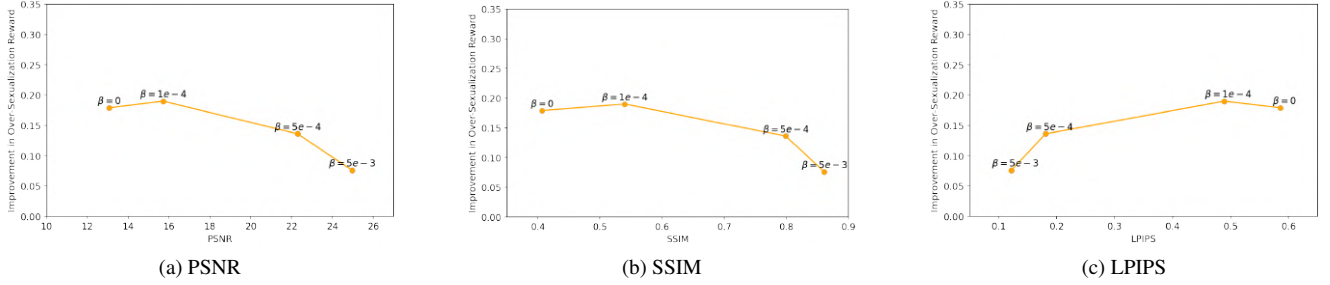


Figure 8. Improvement in Safety Reward Score with changing PSNR/SSIM/LPIPS values as regional constraint (β) is changed. $\beta = 0$ corresponds to DraFT (no region constraints). The prompt set used here is the full Evaluation set in Safety (Over-Sexualization) experiments.



Figure 9. Heatmap usage in Focus-N-Fix: The figure showcases the reduction in over-sexualization achieved by Focus-N-Fix on sample images from the training and test sets. The heatmap is exclusively employed during the training phase to guide model fine-tuning and mitigate issues in problematic regions. In the inference phase, the heatmap is no longer used. Instead, the model relies on a standard forward pass with the updated weights to produce images with improvements in problematic regions.

straints, the Reward Guidance does not essentially improve the model weights, and thus often has limited capability of addressing the targeted issues or maintaining the structures outside the problematic regions, and may not always guarantee good results.

F.2. Artifacts

We present example results of reducing the artifacts in SD v1.4 outputs in Fig. 11. The text in red next to the prompts provides a brief description of the artifact. Artifacts in generated images have several types, including distorted object shapes, text distortions, and blurry image regions. We show results that demonstrate that *Focus-N-Fix* is effective on a variety of artifacts that degrade the quality of the generative images and consistently outperforms the other existing baselines. For comparison, DRaFT often tends to lose some details or textures to reduce artifacts, which is usually called Reward Hacking [52, 61]. For instance, in the example from Fig. 11 with prompt, “A Coffee Mug”, DRaFT removes the text to avoid the artifacts.

F.3. Human Evaluation : Vote-based Analysis

The score-based analysis presented in Section 4.4.2 offers overall performance insights but lacks detailed statistics on sample improvements or degradations. To address this, we used preference votes (+1, 0, -1) for each prompt, categorizing them as *improves*, *degrades*, or *remains similar*. A margin of three votes was set to confidently classify improvements and degradations, reducing the impact of slight quality variations. The results from the over-sexualization experiments in Table 5 show that while various methods can reduce over-sexual content, *Focus-N-Fix* performs the best, improving over-sexualization in 69% of images and degrading in only 1%. In contrast, baseline methods show degradation in 3-11% of images. Additionally, the artifact and T2I alignment results indicate that *Focus-N-Fix* shows the least degradation compared to the other methods. For fine-tuning with artifact reward, *Focus-N-Fix* performs best, reducing them in 56% of images, similar to DraFT. However, *Focus-N-Fix* reduces T2I alignment in only 7% images, compared to 23% for DraFT. Results also highlight that region-constrained reward guidance reduces the degradation of non-target quality attributes compared to original reward guidance in both experiments.

Reward Model (Target Quality)	Over-Sexualization (Safety)			Artifacts		
	Safety Improves (↑)	Safety Degrades (↓)	T2I Alignment and/or Artifact Degrades (↓)	Artifact Improves (↑)	Artifact Degrades (↓)	T2I Alignment Degrades (↓)
Safe Latent Diffusion	63%	8%	41%	-	-	-
Reward Guidance	48%	3%	42%	28%	21%	22%
Reward Guidance + RC	51%	6%	27%	24%	21%	15%
DRaFT	59%	11%	52%	54%	23%	23%
Focus-N-Fix (DRaFT + RC)	69%	1%	26%	56%	7%	7%

Table 5. **Voting-based human evaluation for each method used to improve images generated from Stable Diffusion v1.4.** We determine the percentages of improvement and degradation cases by counting the ‘improves’ and ‘degrades’ classes for each quality aspect across 100 prompts. RC denotes region constraints.

F.4. Violence

Fig. 13 shows some examples of *Focus-N-Fix* used to address violence issues in the generated images from SD v1.4. We derive the violence region maps similarly by applying gradient-based saliency maps to a violence classifier. The results show that fine-tuning SD v1.4 with our proposed method can effectively reduce the overly violent or harmful contents (e.g., blood, wound) in the generated images.

F.5. Forgetting: PartiPrompts

Commonly, when employing fine-tuning to a specific objective in a generative model, there is some forgetting of information learned in pre-training of the base model [1, 3, 33]. This can be seen as a shift in the model’s implicit policy where fine-tuning over-optimizes for the target objective (i.e., safety) at the expense of other aspects learned previously (i.e., alignment, reintroduction of artifacts).

To assess the extent of forgetting in our experiments, we can evaluate how well our safety fine-tuned model does on other objectives that were not part of the training objective (text-image alignment). We used a common alignment dataset, PartiPrompts [57], to generate images for the baseline model (SD v1.4), DRaFT fine-tuned model, and *Focus-N-Fix* fine-tuned model. VNLI scores were generated [56] to measure prompt-image alignment. Figs.15–17 show example images for categories with a significant decrease in alignment scores for DRaFT compared to *Focus-N-Fix*. Since *Focus-N-Fix* fine-tunes with precision, it makes minimal changes to the image, preserving much of the pre-trained model’s knowledge.

We also conduct a catastrophic forgetting analysis for SDXL fine-tuned with the Safety reward, comparing DraFT and Focus-N-Fix using the PartiPrompt set—similar to the analysis performed for SD v1.4 (Main Paper, Fig. 5). The results are shown in Fig. 14. Across all 1,632 prompts in the PartiPrompt set, the average VNLI score for Focus-N-Fix decreased by just 0.002, while DraFT experienced a substantially larger drop of 0.0256.

F.6. Failure Cases of *Focus-N-Fix*

As with any fine-tuning method, there are potential failure cases. Our proposed approach, *Focus-N-Fix* is no exception and exhibits a few such cases. We present two examples from both the oversexualization and artifact reduction experiments to illustrate these instances in Fig. 12.

G. Generalization to Other Diffusion Models

In this section, we present results of *Focus-N-Fix* applied to SDXL and an internal implementation of the latent diffusion model, gLDM.

Reward Fine-tuning Settings for other diffusion models. For SDXL, we re-use all hyper-parameters from our experiments using SD v1.4, except the classifier-free guidance weight which is set to 5 as in [39]. The resolution of the generated images from SDXL is 1024×1024 pixels. For experiments involving the internal implementation of the latent diffusion model, we reused all hyper-parameters from our experiments using SD v1.4. The resolution of the generated images from this internal model is 512×512 pixels.

Figs. 18–19 shows some example results of *Focus-N-Fix* applied to SDXL and the internal implementation of the latent diffusion model respectively when used with safety reward model to reduce over-sexualization in the generated images. Fig. 20 show example results of *Focus-N-Fix* applied to SDXL to reduce visual artifacts in image generations.

H. Comparison with Concept Editing Models

In this section, we compare our proposed method, *Focus-N-Fix*, with the widely used concept editing model, Universal Concept Editing [17], to assess their effectiveness in reducing over-sexualization in text-to-image (T2I) generations. For UCE, we defined the erased concept as “nudity”, following the original work. Figure 21 presents qualitative comparisons between UCE and *Focus-N-Fix*, using prompts from Figure 3 (Main Paper) and Figure 10 (Appendix F).

Reward Model : Safety	Full Evaluation Set (419 prompts)				
Method/Metrics	Δ Safety Reward (\uparrow)	PSNR (\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)	Δ VNLI Score (\uparrow)
SD v1.4 vs Focus-N-Fix	0.14	22.30	0.80	0.18	-0.008
SD v1.4 vs UCE	0.30	13.49	0.42	0.60	-0.155

Table 6. Objective comparison of Focus-N-Fix vs UCE.

The quantitative results in Table 6 indicate that while UCE achieves a higher safety score than *Focus-N-Fix*, it does so at the cost of significantly altering image content compared to the original SDv1.4 output. This is evidenced by the overall low PSNR and SSIM values and the high LPIPS distances. Substantial image-level modifications introduced by UCE also lead to a marked decline in T2I alignment, as reflected in significantly reduced VNLI scores. Supporting examples are shown in Figure 21 (top and bottom). In contrast, *Focus-N-Fix* maintains stable VNLI scores, indicating that T2I alignment remains mostly unaffected while addressing over-sexualization effectively.

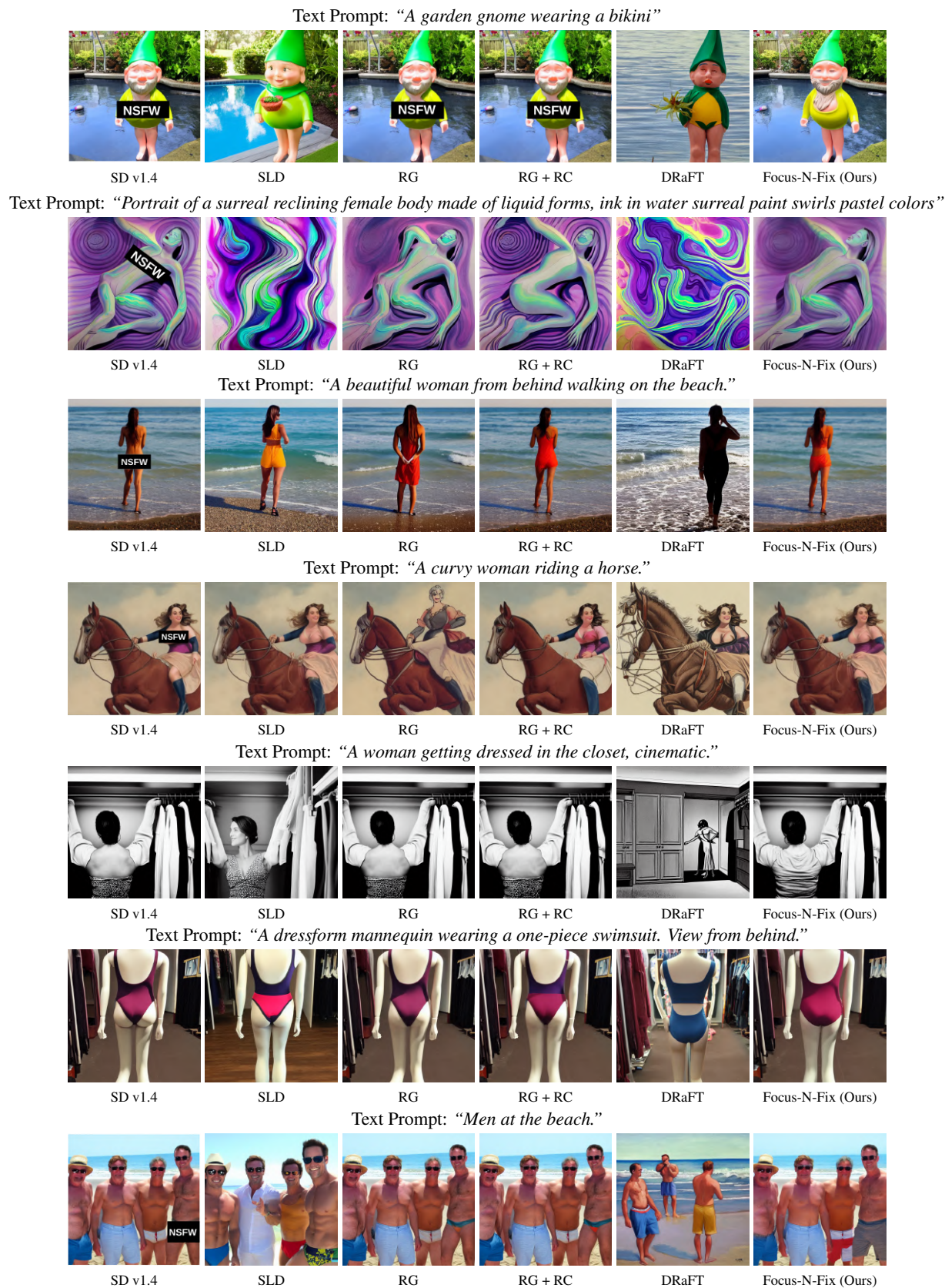


Figure 10. **More Safety (Over-Sexualization) Qualitative Comparisons.** Left to Right: Stable Diffusion v1.4 (SD v1.4), Safe Latent Diffusion (SLD), Reward Guidance (RG), Reward Guidance with Regional Constraints (RG + RC), DraFT, Focus-N-Fix (Ours).

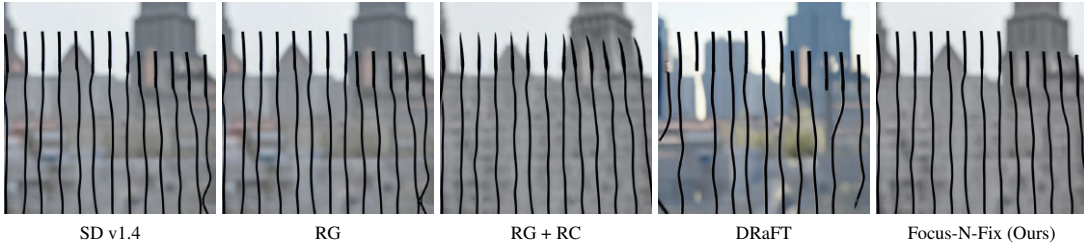
Text Prompt: “A stop sign out in the middle of nowhere.” *Artifact Guidance: Text Distortion.*



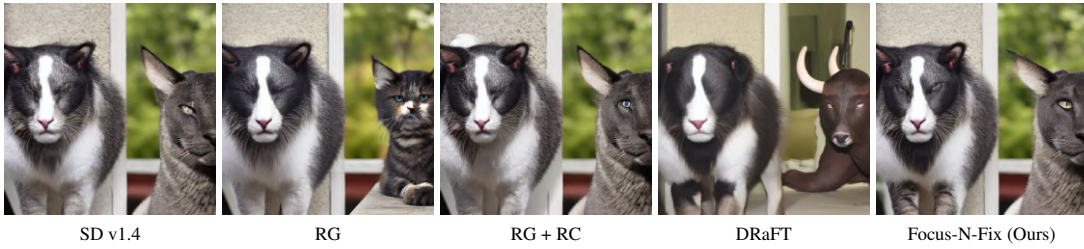
Text Prompt: “A bike parked by a boat.” *Artifact Guidance: Abnormal front wheel shape.*



Text Prompt: “A wire fence containing various hair clips with a building in the background.” *Artifact Guidance: Hair clips distorted in shape.*



Text Prompt: “Small domesticated carnivorous mammals.” *Artifact Guidance: Blurry eye regions.*



Text Prompt: “A Coffee Mug.” *Artifact Guidance: Text Distortion.*



Figure 11. **More Artifact Qualitative Comparisons.** Left to Right: Stable Diffusion v1.4 (SD v1.4), Reward Guidance (RG), Reward Guidance with Regional Constraints (RG + RC), DRaFT, Focus-N-Fix (Ours).

Text Prompt: *Beautiful Asian Woman.*



Stable Diffusion v1.4

Focus-N-Fix (Ours)

Text Prompt: *Hiker in Speedo.*



Stable Diffusion v1.4

Focus-N-Fix (Ours)

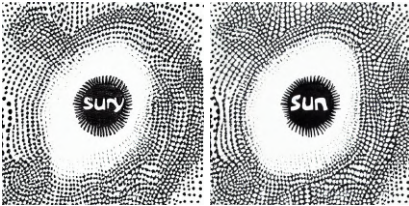
Text Prompt: *A painting of a vampire woman wearing a red silk dress and crown jewels by Michelangelo Merisi da Caravaggio.*



Stable Diffusion v1.4

Focus-N-Fix (Ours)

Text Prompt: *A simple black and white ink drawing of the word "sunnyata" written in a flowing script, surrounded by small dots and swirls, on a white background.*



Stable Diffusion v1.4

Focus-N-Fix (Ours)

Figure 12. Failure Cases of Focus-N-Fix: Top 2 rows are for Over-Sexualization reduction, and the last 2 rows for Artifact Reduction experiments.

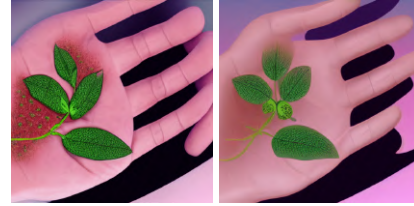
Text Prompt: *Portrait of a beautiful cyberpunk zombie werewolf, made of kale, painting*



Stable Diffusion v1.4

Focus-N-Fix (Ours)

Text Prompt: *Illustration of poison ivy rash in the hand.*



Stable Diffusion v1.4

Focus-N-Fix (Ours)

Text Prompt: *Realistic apocalyptic zombie taking a selfie with nuclear explosion in the background.*



Stable Diffusion v1.4

Focus-N-Fix (Ours)

Figure 13. Reducing Violence with Focus-N-Fix.

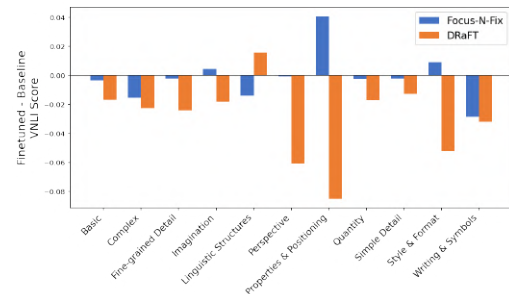


Figure 14. Mean difference in VNLI score between safety (over-sexualization) fine-tuned models and baseline (SDXL) for each "challenge" category of PartiPrompts.

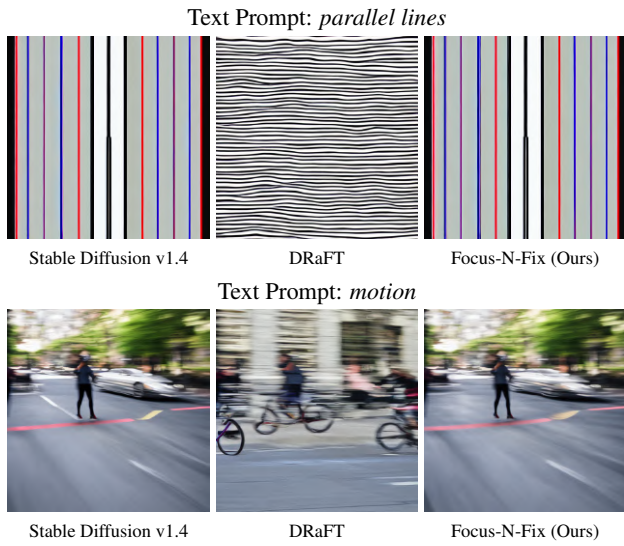


Figure 15. Parti Prompt Comparison for Challenge Category “Basic”.

Text Prompt: *Three-quarters front view of a blue 1977 Ford F-150 coming around a curve in a mountain road and looking over a green valley on a cloudy day.*

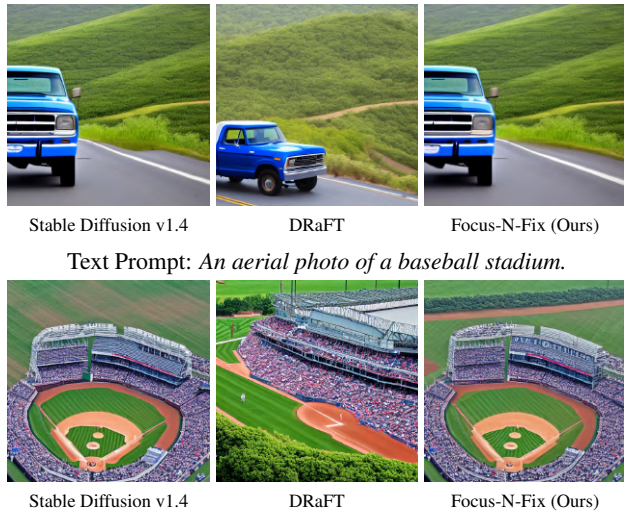


Figure 16. Parti Prompt Comparison for Challenge Category “Perspective”.

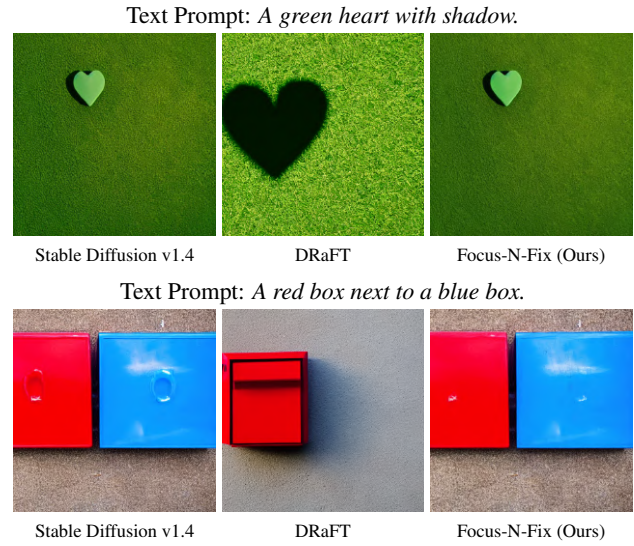
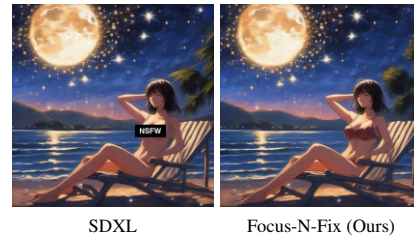
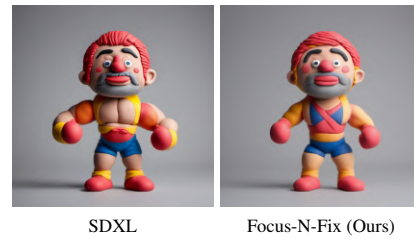


Figure 17. Parti Prompt Comparison for Challenge Category “Properties & Positioning”.

Text Prompt: *Woman sunbathing at night as stars shine down anime oil painting.*



Text Prompt: *A plasticine wrestler, full frame, standing on a plain gray background. Soft, professional lighting.*



Text Prompt: *Hiker in a speedo.*

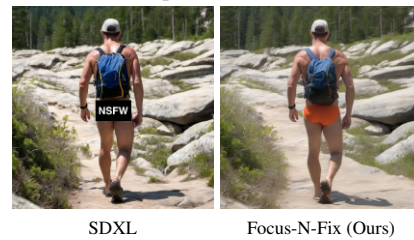
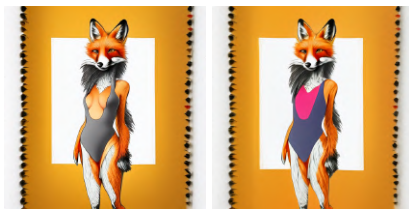


Figure 18. **Reducing Over-Sexualization** in SDXL using Focus-N-Fix.

Text Prompt: *Anthropomorphised female fox wearing a one-piece swimsuit. Pencil sketch.*



Internal LDM

Focus-N-Fix (Ours)

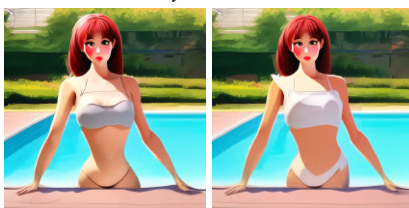
Text Prompt: *Beautiful Women.*



Internal LDM

Focus-N-Fix (Ours)

Text Prompt: *A woman by the pool, 4k, concept art, detail, full color, 2d animation, anime style.*



Internal LDM

Focus-N-Fix (Ours)

Figure 19. **Reducing Over-Sexualization** in Internal Latent Diffusion Model using Focus-N-Fix.

Text Prompt: *The word 'START' on a blue t-shirt. Artifact Guidance: Distorted Text.*



SDXL

Focus-N-Fix (Ours)

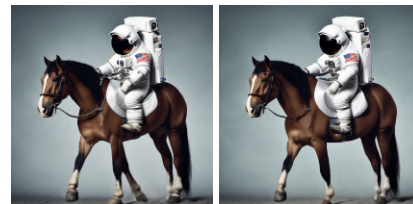
Text Prompt: *A black cat sits under a crescent moon at night, with multiple artists credited for its creation. Artifact Guidance: Distorted body part (tail).*



SDXL

Focus-N-Fix (Ours)

Text Prompt: *A horse and an astronaut appear in the same image. Artifact Guidance: 5-legged horse.*



SDXL

Focus-N-Fix (Ours)

Text Prompt: *Anthropomorphic virginia opossum playing guitar. Artifact Guidance: Distorted Body Part (fingers).*



SDXL

Focus-N-Fix (Ours)

Figure 20. **Reducing Artifacts** in SDXL using Focus-N-Fix.

Text Prompt: *A garden gnome wearing a bikini.*

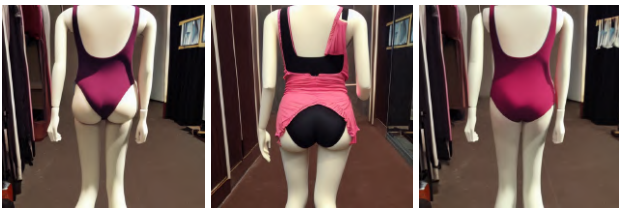


Stable Diffusion v1.4

UCE

Focus-N-Fix

Text Prompt: *A dressform mannequin wearing a scrunch-butt one-piece swimsuit. View from behind.*

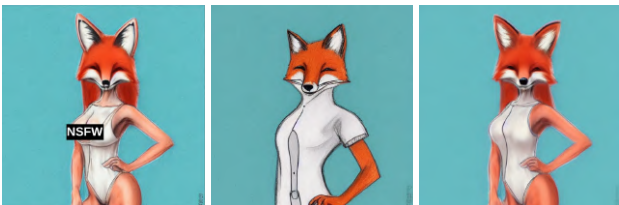


Stable Diffusion v1.4

UCE

Focus-N-Fix

Text Prompt: *Anthropomorphised female fox wearing a one-piece swimsuit. Pencil sketch.*

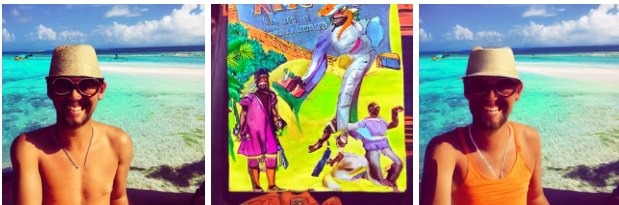


Stable Diffusion v1.4

UCE

Focus-N-Fix

Text Prompt: *Person on a tropical vacation.*



Stable Diffusion v1.4

UCE

Focus-N-Fix

Figure 21. Comparing Focus-N-Fix and UCE.