

DIFFUSION-NPO: NEGATIVE PREFERENCE OPTIMIZATION FOR BETTER PREFERENCE ALIGNED GENERATION OF DIFFUSION MODELS

Anonymous authors

Paper under double-blind review

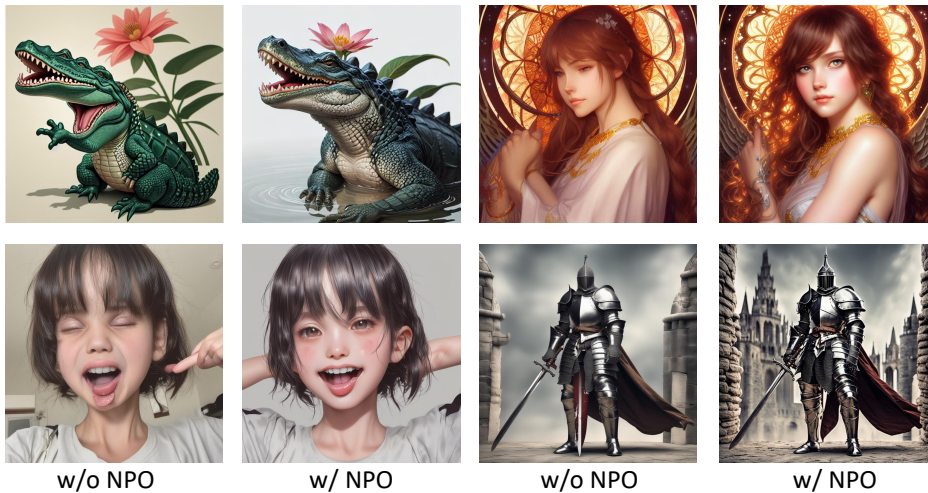


Figure 1: Diffusion-NPO enhances high-frequency details, color and lighting, and low-frequency structures in images by aligning human’s negative preference.

ABSTRACT

Diffusion models have made substantial advances in image generation, yet models trained on large, unfiltered datasets often yield outputs misaligned with human preferences. Numerous methods have been proposed to fine-tune pre-trained diffusion models, achieving notable improvements in aligning generated outputs with human preferences. However, we argue that existing preference alignment methods neglect the critical role of handling unconditional/negative-conditional outputs, leading to a diminished capacity to avoid generating undesirable outcomes. This oversight limits the efficacy of classifier-free guidance (CFG), which relies on the contrast between conditional generation and unconditional/negative-conditional generation to optimize output quality. In response, we propose a straightforward but versatile effective approach that involves training a model specifically attuned to negative preferences. This method does not require new training strategies or datasets but rather involves minor modifications to existing techniques. Our approach integrates seamlessly with models such as SD1.5, SDXL, video diffusion models and models that have undergone preference optimization, consistently enhancing their alignment with human preferences.

1 INTRODUCTION

Diffusion models have made significant strides in image/video generation (Rombach et al., 2022; Podell et al., 2023; Dhariwal & Nichol, 2021; Singer et al., 2022; Shi et al., 2024; Wang et al., 2024; Blattmann et al., 2023; Liang et al., 2024a). However, diffusion models trained on massive unfiltered image-text pairs (Schuhmann, 2022; Sun et al., 2024) often generate results that do not

054 align with human preferences. To address this issue, many methods (Wu et al., 2023; 2024) have
055 been proposed to align diffusion models with human preferences, aiming to drive the generation to
056 better match what users desire.

057 Human preference alignment methods typically require the prior collection of a human preference
058 dataset, such as Pick-a-pic (Kirstain et al., 2023). The standard procedure involves gathering pairs
059 of images generated from the same prompt and annotating them according to human preferences.
060 Rather than assigning direct scores, these preferences are usually ranked in order. This ranking
061 is then utilized to train a scoring/reward model for text-image pairs using a contrastive loss func-
062 tion (Ouyang et al., 2022). To explore this topic in depth, we first review existing approaches for
063 aligning diffusion models with human preferences. In general, current methods can be categorized
064 into three types:

- 065
- 066 a) **Differentiable Reward (DR)**: These approaches directly feed multi-step generated images
067 into a pretrained reward model, updating the diffusion models through gradient backpropa-
068 gation (Xu et al., 2024; Prabhudesai et al., 2024; Zhang et al., 2024; Wu et al., 2023; 2024).
069 While simple and direct, these methods are prone to reward leakage (Zhang et al., 2024).
 - 070 b) **Reinforcement Learning (RL)**: In these approaches, the denoising process of diffusion
071 models is formulated as an equivalent Markov decision process (MDP) (Puterman, 2014).
072 PPO (Schulman et al., 2017) and its variants are typically adopted for preference optimiza-
073 tion. Images are generated and evaluated online based on the reward feedback, aiming to
074 increase the probability of generating high-reward images. These approaches employ SDE
075 solvers to achieve stochastic sampling and importance sampling (Sutton, 2018).
 - 076 c) **Direct Preference Optimization (DPO)**: These approaches simplify the reinforc-
077 e-ment learning training objective into a straightforward simulation-free training objec-
078 tive (Rafailov et al., 2024; Wallace et al., 2024). They do not require training reward mod-
079 els, nor do they need online generation and sampling; instead, they only require fine-tuning
080 on pre-collected paired preference data. Although simple, these approaches often under-
081 perform reinforcement learning-based methods, especially for out-of-distribution inputs.

082

083 Despite previous efforts to make models generate human-aligned images, we raise an important
084 question: *How can a model know to avoid generating poor images if it only knows how to generate*
085 *good ones without understanding what is bad?*

086 We identify a crucial oversight in current diffusion model preference alignment efforts: most dif-
087 fusion generation rely heavily on the classifier-free guidance (CFG) (Ho & Salimans, 2022; Karras
088 et al., 2024; Shen et al., 2024; Ahn et al., 2024). CFG requires the model to simultaneously com-
089 pute outputs under both conditional inputs and negative-conditional/unconditional inputs at each
090 denoising step, then linearly combine these outputs to bias the final prediction towards the condi-
091 tional inputs and away from the negative-conditional inputs. Ideally, we expect the model’s output
092 under the conditional inputs to align closely with human preferences, while the output under the
093 negative-conditional inputs should diverge from human preferences to maximize preference align-
094 ment. However, previous works focus exclusively on training models to generate outputs that align
095 with human preferences, without considering the equally important task of teaching models to recog-
096 nize and avoid generating outputs that humans do not favor. This oversight limits the effectiveness
097 of existing alignment strategies, particularly in scenarios where distinguishing between preferred
098 and non-preferred outputs is crucial.

099 To address this issue, we propose **Negative Preference Optimization (NPO)**: training an additional
100 model that is aligned with preferences opposite to human. Importantly, our crucial insight is that
101 *training such a negative preference aligned model requires no new training strategies or datasets,*
102 *only minor modifications to existing methods.* 1) Approaches like differential reward and reinforc-
103 e-ment learning, all need a reward model for training. We simply multiply the output of reward model
104 by -1 , which allows us to train a negative preference model using the same approaches. 2) For
105 DPO-based methods, we reverse the order of the preferred image pairs. Notably, during the training
106 of the reward model applied for differential reward and reinforcement learning approaches, the im-
107 age order can also be reversed to train the reward model. Therefore, in essence, all strategies can be
perceived as reversing the order of image pairs in the collected preference data by adapting the same
training procedure. Fig. 2 provides an overview of our method.

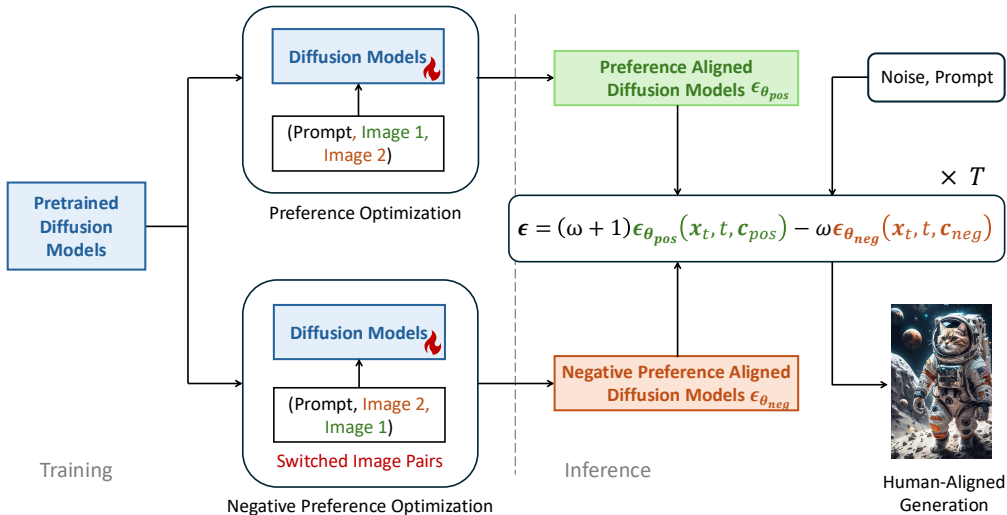


Figure 2: High-level overview of negative preference optimization (NPO). (Training) NPO needs no new training strategies and datasets. NPO training can be achieved through switching preference image pairs with existing preference optimization methods. (Inference) NPO trained models serve as the unconditional/negative-conditional predictors in the classifier-free guidance.

We validate the effectiveness of NPO on text-to-image generation with SD1.5 (Rombach et al., 2022) and SDXL (Podell et al., 2023) and text-to-video generation with VideoCrafter2 (Chen et al., 2024). Our model can be used in a plug-and-play manner with these baseline models and their various preference-optimized versions, consistently improving generation quality. Fig. 3 shows our comparative results. We evaluate our method using the widely adopted Pick-a-pic validation set, scoring with metrics including HPSv2, ImageReward, PickScore, and LAION-Aesthetic. Our approach significantly improves performance across all metrics.

2 UNDERSTANDING CLASSIFIER-FREE GUIDANCE

Preliminary of CFG. CFG has become a necessary and important technique for improving generation quality and text alignment of diffusion models. For convenience, we focus our discussion on the general formal of diffusion models, *i.e.*, $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon$ (Kingma et al., 2021). Suppose we learn a score estimator from a epsilon prediction neural network $\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}, t)$, and we have $\nabla_{\mathbf{x}_t} \log \mathbb{P}_{\theta}(\mathbf{x}_t | \mathbf{c}; t) = -\frac{\epsilon_{\theta}(\mathbf{x}_t, t)}{\sigma_t}$. The sample prediction at timestep t of the score estimator is formulated as

$$\hat{\mathbf{x}}_0 = \frac{1}{\alpha_t} (\mathbf{x}_t + \sigma_t^2 \nabla_{\mathbf{x}_t} \log \mathbb{P}_{\theta}(\mathbf{x}_t | \mathbf{c}; t)). \quad (1)$$

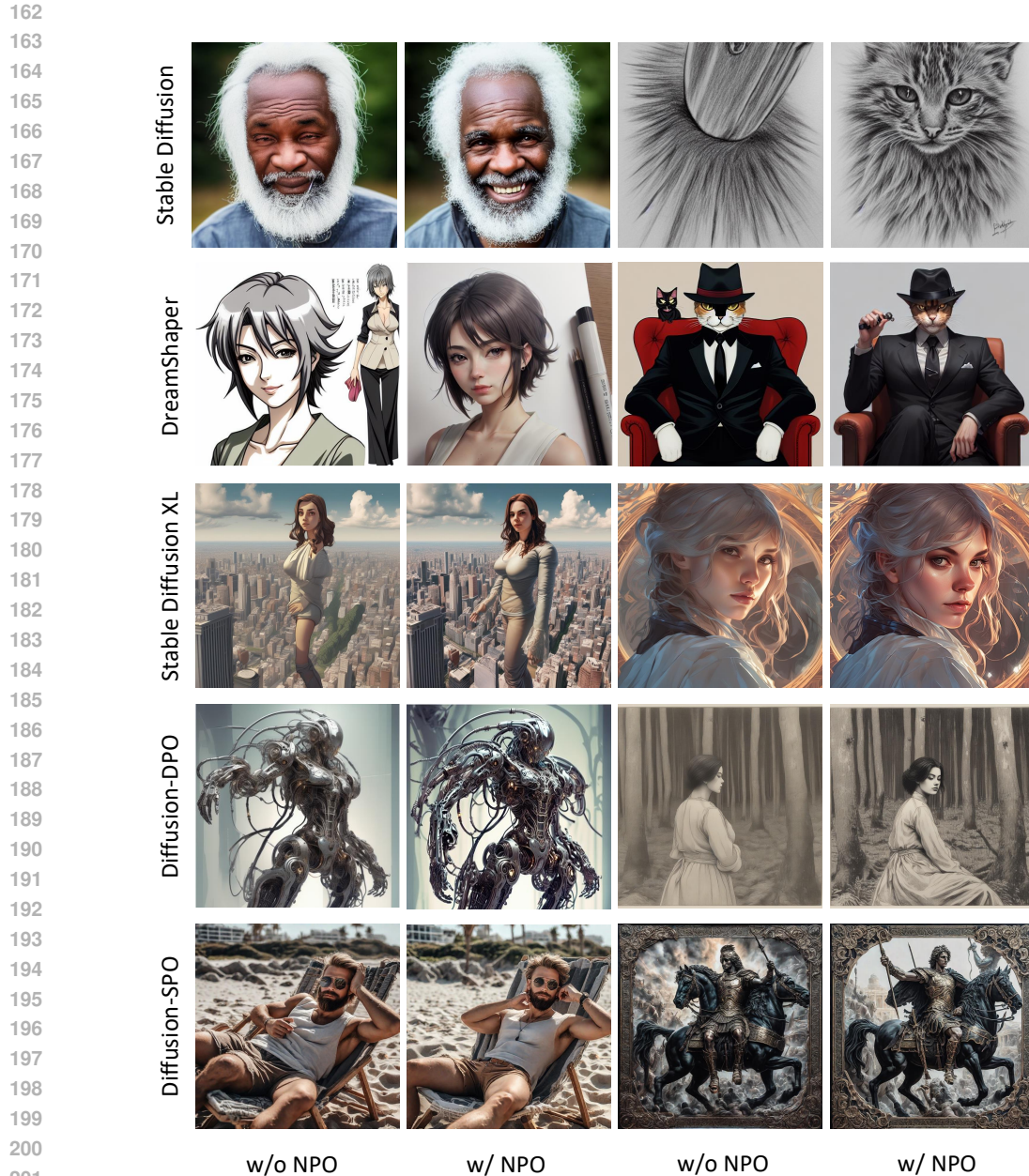
Applying the CFG is equivalent to add an additional score term (Karras et al., 2024), that is, we replace $\nabla_{\mathbf{x}_t} \log \mathbb{P}_{\theta}(\mathbf{x}_t | \mathbf{c}; t)$ in Eq. 1 with the following term,

$$\nabla_{\mathbf{x}_t} \log \mathbb{P}_{\theta}(\mathbf{x}_t | \mathbf{c}; t) + \nabla_{\mathbf{x}_t} \log \left[\frac{\mathbb{P}_{\theta}(\mathbf{x}_t | \mathbf{c}; t)}{\mathbb{P}_{\theta}(\mathbf{x}_t | \mathbf{c}'; t)} \right]^{\omega}, \quad (2)$$

where ω is to control the strength of CFG, \mathbf{c} and \mathbf{c}' are conditional and unconditional/negative-conditional inputs, respectively. It is apparent that the generation will be pushed to high probability region of $\mathbb{P}_{\theta}(\mathbf{x}_t | \mathbf{c}; t)$ and relatively low probability region of $\mathbb{P}_{\theta}(\mathbf{x}_t | \mathbf{c}'; t)$. Write the above equation into the epsilon format, and then we have

$$\epsilon_{\theta}^{\omega} = (\omega + 1)\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}, t) - \omega\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}', t). \quad (3)$$

Motivating example. To maximize human preference alignment, in Eq. 3, the green component should guide the generated results to closely match human preferences, while the orange component



203 Figure 3: NPO works as a plug-and-play inference enhancement strategy. It can be easily combined
204 with base diffusion models and preference optimized diffusion models for better human preference
205 aligned generation. Zoom out for better comparison in details.

206
207
208
209 should direct the results away from undesired outcomes. However, most preference optimization
210 methods focus exclusively on optimizing the green component, neglecting the orange component
211 and thereby weakening its impact. To illustrate this point, we setup a motivating experiment to
212 investigate the influence of the orange component. We employ two baselines:

- 213
214
215
1. We use the DPO-optimized SD1.5 (Wallace et al., 2024; Rombach et al., 2022) for both the green component and the orange component.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

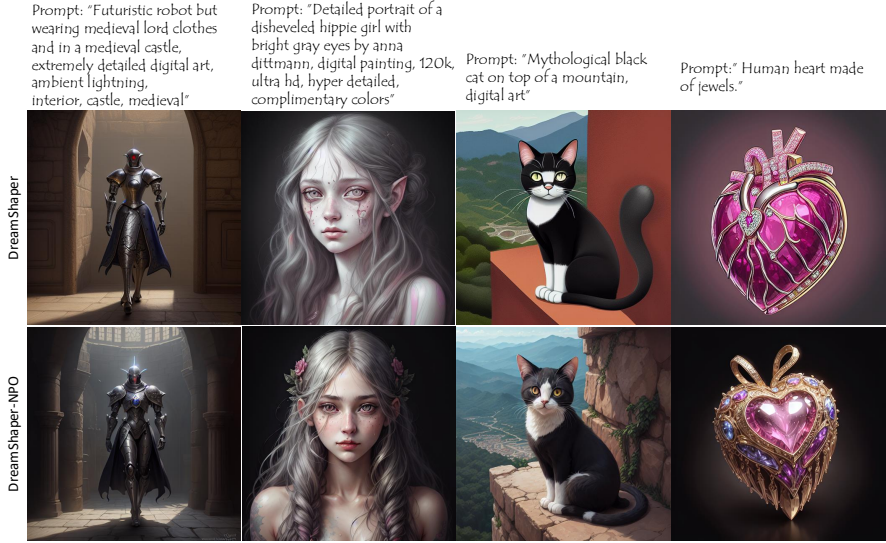


Figure 4: Plug-and-play NPO on DreamShaper. NPO not only works on the base Stable Diffusion and its preference optimized variants, but also works on improving customized model finetuned on high-quality data.

2. We use the DPO-optimized SD1.5 for the green component and the model merged from weights of DPO-optimized SD1.5 (0.6×) and original SD1.5 (0.4×) for the orange component, generating results with the same seed.

We compare the generated images from the two baselines one by one, score them using HPSv2, ImageReward, PickScore, and LAION-Aesthetic, and calculate their win probabilities. The results are shown in Fig. 5. We can observe a significant improvement in human preference compared to only using the DPO-optimized model.

Analysis: the weight merge is an approximated NPO. What is the meaning of the weight merged model? Suppose the weight of original SD is θ , and then the DPO-optimized model weight can be denoted as $\theta + \eta$ since it is further trained from the original weight θ . The merged model weight is $\gamma(\theta + \eta) + (1 - \gamma)\theta = \theta + \gamma\eta = \theta + \eta + (1 - \gamma)(-\eta)$, (4)

where $\gamma \in [0, 1]$ is the merge factor. We can observe that after weight fusion, the weight offset obtained through DPO optimization η has decreased. Consequently, the DPO weight offset η has a weaker impact on the generated results, enabling the model to output results that are more contrary to human preferences. Replace $(1 - \gamma)(-\eta)$ with δ , and then the weight can be represented as

$$\theta' = \theta + \eta + \delta. \quad (5)$$

The above equation decomposes the weight applied for negative-conditional predictions into three parts: the original model weight θ , the preference alignment weight offset (direction) η , a weight offset opposite to the preference alignment δ . Our paper aims to train a suitable δ and investigate its properties. Note that, once the η and δ are obtained, we can also flexibly change the influence of each weight offset by multiply simple scale factors. That is,

$$\theta' = \theta + \alpha\eta + \beta\delta, \quad (6)$$

where $\alpha, \beta \in [0, 1]$.

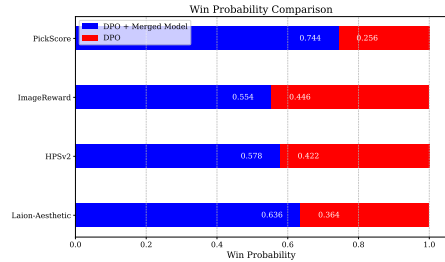


Figure 5: Motivating example results of Section 2. Applying merged model as the orange component (*i.e.*, prediction for unconditional/negative-conditional inputs) effectively improves the human preference alignment.

Table 1: Quantitative performance comparison with stable diffusion v1-5 based models. * means the metrics are copied from SPO papers. Other metrics are tested with official weights.

Method	PickScore	HPSv2	ImageReward	Aesthetic
SD-1.5	20.75	26.84	0.1064	5.539
*DDPO	21.06	24.91	0.0817	5.591
*D3PO	20.76	23.97	-0.1235	5.527
Diff.-DPO	20.98	25.05	0.1115	5.505
Diff.-SPO	21.41	26.85	0.1738	5.946
SD-1.5 + NPO	21.26	27.36	0.2028	5.667
Diff.-SPO + NPO	21.65	27.09	0.1939	5.999
Diff.-DPO + NPO (reg= 500)	21.58	27.60	0.3101	5.762
Diff.-DPO + NPO (reg= 1000)	21.43	27.36	0.3472	5.773
DreamShaper	21.96	27.97	0.7131	6.085
DreamShaper + NPO ($\alpha = 1.0$)	22.38	28.31	0.7396	6.169
DreamShaper + NPO ($\alpha = 0.6$)	22.46	28.08	0.6626	6.496

3 NEGATIVE PREFERENCE OPTIMIZATION

Previous approaches primarily focus on training single model weight that aligns with human preferences. However, these methods often overlook the importance of unconditional outputs of classifier-free guidance in the diffusion generation process. Our approach seeks to train a weight offset δ that opposes human preferences to fulfill the role of unconditional outputs. By integrating this offset with the base model’s weights, it functions as a predictor for unconditional inputs, thereby reducing the likelihood of generating outputs that conflict with human preference. The important motivation for negative preference optimization is that a preference-aligned model should not only learn to generate desirable outcomes but also understand what constitutes undesirable ones. This dual understanding is crucial for maximizing preference alignment while minimizing the occurrence of unwanted results.

3.1 TRAINING WITH NPO

An important insight in our work is that achieving negative preference optimization does not require new datasets, reward models, or even new training strategies. Standard preference optimization methods can be directly applied to negative preference optimization.

For methods based on reinforcement learning and differential rewards, which typically rely on a reward model $R(\mathbf{x}, \mathbf{c}) \in [0, 1]$ (can be easily scaled if not in this interval). This reward model can be transformed into the form required for negative preference optimization as follows:

$$R_{\text{NPO}}(\mathbf{x}, \mathbf{c}) = 1 - R(\mathbf{x}, \mathbf{c}). \quad (7)$$

For methods that utilize reward models, we can simply substitute the original $R(\mathbf{x}, \mathbf{c})$ in the algorithm with $R_{\text{NPO}}(\mathbf{x}, \mathbf{c})$.

For methods that train on preference pairs $r = (\mathbf{x}_0, \mathbf{x}_1, \mathbf{c})$, where \mathbf{x}_0 is less preferred and \mathbf{x}_1 is more preferred by humans, and \mathbf{c} is the conditional information used for generation (indicating both images are generated from the same \mathbf{c}), converting this to a negative preference optimization algorithm requires simply reversing the order of the preference pair:

$$r_{\text{NPO}} = (\mathbf{x}_1, \mathbf{x}_0, \mathbf{c}). \quad (8)$$

Beyond the fundamental implementation of negative preference optimization outlined above, it is important to recognize that many preference optimization methods may use CFG during training for sample collection, probability calculation, and gradient backpropagation. NPO can naturally extend to these methods as well. Although these methods might apply CFG during training to bridge the gap between training and inference, they train only a single weight offset, overlooking the fact that the conditional and unconditional (or negative-conditional) outputs in CFG have different optimization objectives (i.e., preference-aligned and negative preference-aligned). This could result in a weight offset that is a compromise between the two opposite objectives, failing to fully achieve preference alignment. We propose to optimize two distinct weight offsets simultaneously.

Table 2: Quantitative performance comparison with stable diffusion XL based models. All metrics are tested with official weights.

Method	PickScore	HPSv2	ImageReward	Aesthetic
SDXL	22.06	27.89	0.6246	6.114
Diff.-DPO	22.57	28.58	0.8767	6.099
Diff.-SPO	22.97	28.58	1.032	6.348
SDXL + NPO	22.32	28.11	0.6831	6.136
Diff.-DPO + NPO	22.69	28.78	0.9210	6.112
Diff.-SPO + NPO	23.08	28.70	1.047	6.438

3.2 INFERENCE WITH NPO

Let θ denote the base model weight, η the weight offset after preference optimization, and δ the weight offset after negative preference optimization. A straightforward strategy is to define $\theta_{pos} = \theta + \eta$ and $\theta_{neg} = \theta + \delta$, and then apply classifier-free guidance as follows:

$$\epsilon_{\theta}^{\omega} = (\omega + 1)\epsilon_{\theta_{pos}}(\mathbf{x}_t, \mathbf{c}, t) - \omega\epsilon_{\theta_{neg}}(\mathbf{x}_t, \mathbf{c}', t). \quad (9)$$

However, this approach often results in a significant output discrepancy between θ_{pos} and θ_{neg} . The outputs from classifier-free guidance should maintain a necessary level of correlation; for example, if two Gaussian noises are completely independent, the variance from the operation above would change from 1 to $2\omega^2 + 2\omega + 1$. We find that it is typically necessary to incorporate the positive weight offset into the negative weights, such that:

$$\theta_{neg} = \theta + \alpha\eta + \beta\delta, \quad \alpha, \beta \in [0, 1] \quad (10)$$

which aligns with our earlier motivating example and analysis.

4 EXPERIMENTS

4.1 VALIDATION SETUP

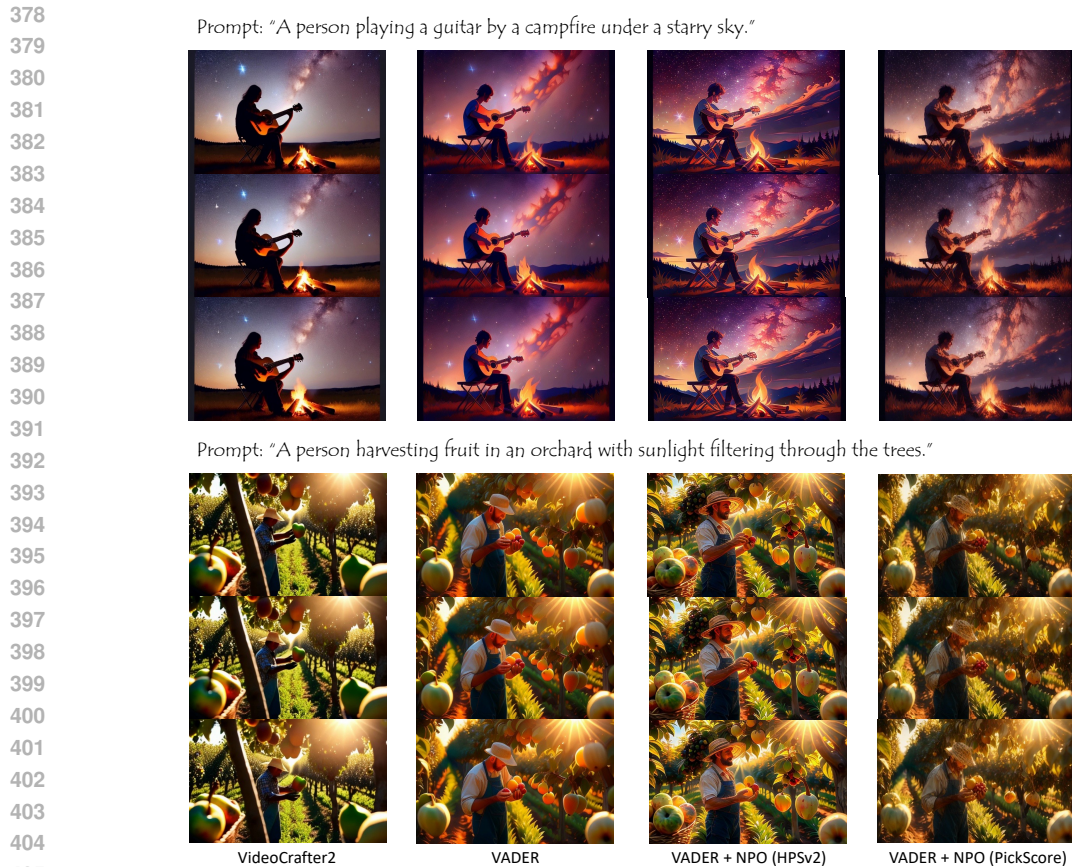
To validate the effectiveness and versatility of our approach, we test it on three baseline methods:

- a) **Diffusion-DPO.** Diffusion-DPO (Wallace et al., 2024) is the first method to incorporate the Direct Preference Optimization (DPO) approach into diffusion training. It introduces a simulation-free and reward model-free training strategy that enables direct training with preference pairs. The effectiveness of this method has been validated on popular text-to-image models, such as the 0.9B Stable Diffusion v1-5 and the 3B Stable Diffusion XL.
- b) **Diffusion-SPO.** Diffusion-SPO (Liang et al., 2024b) combines the DPO approach with reinforcement learning. It involves online sample generation, stochastic solvers, and probability calculations, while utilizing the DPO optimization objective for training. This method requires a reward model to score preferences for generated images online. Its effectiveness has also been demonstrated on the 0.9B Stable Diffusion v1-5 and the 3B Stable Diffusion XL for text-to-image generation.
- c) **VADER.** VADER (Prabhudesai et al., 2024) is a differential reward-based approach that has shown effectiveness in text-to-video generation, significantly enhancing the aesthetic quality of generated videos from raw models.

Therefore, our validation baselines include differential reward, reinforcement learning, and direct preference optimization (the three typical kinds of methods we mentioned), covering both text-to-image and text-to-video tasks. We believe our validation is sufficiently convincing to demonstrate the effectiveness of our approach. Unless otherwise specified, we use the default training and inference configurations for all the aforementioned methods, including training data, number of training iterations, CFG strength, etc.

4.2 COMPARISON

Quantitative comparison. For text-to-image generation, we conduct the quantitative evaluation of our method by following previous work and using the Pick-a-pic ‘test_unique’ split as the test-



406 Figure 6: Video comparison. The videos are trained using 12 frames. For better visualization, we
 407 sample one key frame from every four frames..
 408

409
 410 ing benchmark (Kirstain et al., 2023). We employ PickScore (Kirstain et al., 2023), HPSv2 (Wu
 411 et al., 2023), ImageReward (Xu et al., 2024), and Laion-Aesthetic (Schuhmann, 2022) as evalu-
 412 ation metrics. The results of the quantitative evaluation are summarized in Tables 1 and 2. The
 413 tables demonstrate that NPO, when combined with the base model and its preference-optimized ver-
 414 sions, consistently enhances the aesthetic quality of the generated results. In addition to reporting
 415 the average scores, as illustrated in Fig. 7, we calculate the proportion of samples generated with
 416 the same prompt that achieve a higher preference score. The results generated using NPO signifi-
 417 cantly outperform those without NPO. For text-to-video generation, we compare four baselines:
 418 VideoCrafter2, VADER, VADER + NPO (HPSv2), and VADER + NPO (PickScore). Among these,
 419 VADER + NPO (HPSv2) is optimized using both HPSv2 and Laion-Aesthetic as reward models,
 420 while VADER + NPO (PickScore) is optimized using PickScore as the reward model. We train the
 421 models using animal-related prompts, as was done with VADER, and evaluate on unseen animal-
 422 related prompts (same domain) and additional human prompts (out domain). The results, presented
 423 in Table 3, reveal that VADER + NPO (HPSv2) shows significant improvements across all four met-
 424 rics, particularly in the HPS and Laion-Aesthetic metrics. VADER + NPO (PickScore) demonstrates
 425 greater improvement in the PickScore metric and, on animal-related prompts, even achieves better
 426 HPSv2 performance than VADER + NPO (HPSv2).

426 **Qualitative comparison.** Fig. 3, Fig. 4, Fig. 6, Fig. 11, Fig. 12, Fig. 13, Fig. 14 and Fig. 15, present
 427 a comparison of results generated with and without NPO across various scenarios. We observe that
 428 NPO significantly enhances high-frequency details, color and lighting, and low-frequency structures
 429 in images, consistently improving human preference scores.

430 **User preference.** We assess the generation quality in three specific areas: Color and Lighting,
 431 High-Frequency Details, and Low-Frequency Composition. For Color and Lighting, users evaluate
 whether the generated images display natural and visually pleasing color schemes and lightings. For

Table 3: Quantitative performance comparison on text-to-video generation. All metrics are tested with official weights. Avg means the average score. Win means the average winning ratio to other methods. HPSv2 means we apply both aesthetic predictor and HPSv2 for training. PickScore means we apply PickScore for training.

Method	Aesthetic		HPSv2		ImageReward		PickScore	
	Avg	Win	Avg	Win	Avg	Win	Avg	Win
Animal								
VideoCrafter2	5.527	0.00%	29.65	2.08%	1.368	30.73%	22.44	16.81 %
VADER	6.154	55.21%	32.24	46.88%	1.486	58.33%	22.97	34.23%
VADER + NPO (PickScore)	6.110	50.00%	32.81	82.81%	1.463	53.65%	24.16	98.44%
VADER + NPO (HPSv2)	6.379	94.79%	32.52	68.23%	1.492	59.96%	23.14	50.52%
Human								
VideoCrafter2	5.726	1.04%	27.92	10.27%	0.9583	33.33%	22.41	27.75%
VADER	6.462	61.46%	29.74	51.71%	1.102	46.35%	22.55	34.23%
VADER + NPO (PickScore)	6.244	39.58%	29.58	51.71%	1.086	55.73%	23.35	89.06%
VADER + NPO (HPSv2)	6.855	97.92%	30.76	86.98%	1.164	64.58%	22.71	48.29%

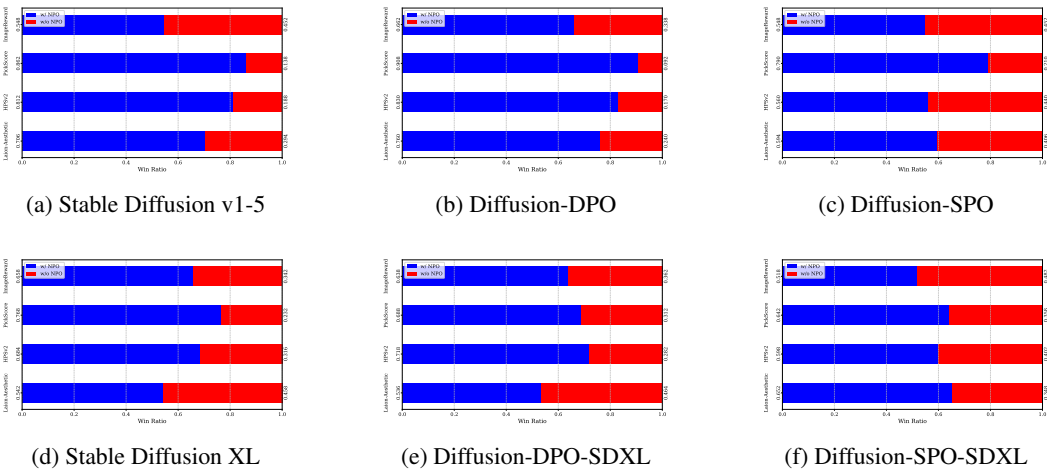


Figure 7: Quantitative winning ratios.

High-Frequency Details, users assess the level of detail in textures and the sharpness of fine features, such as edges and small-scale elements. For Low-Frequency Composition, users examine the overall structure and balance of the images. We conduct the user study using the prompts from Pickcapic ‘validation_unique’, with different models generating images based on the same random seed. Users evaluate the models on the three aspects mentioned above and have three choices: “No Preference” (draw), “NPO is better,” or “NPO is worse.” We distribute questionnaires to 15 volunteers online, with each questionnaire containing 50 pairs of generated images (randomly sampled from SDXL, SD15, and DreamShaper). A total of 750 votes are collected. The final results are presented in the Fig. 8. The user study indicates that NPO significantly enhances high-frequency details in the generated outputs, while also producing colors and lighting that align more closely with human preferences. Additionally, NPO can improve the compositional structure of the generated images to some extent.

Hyper-parameter sensitivity analysis. Negative preference optimization involves a crucial trade-off regarding how much the unconditional/negative-conditional outputs deviate from the conditional outputs. If the deviation is too small, the optimization becomes ineffective; if too large, it may result in blurred or unnatural images. During training, this trade-off is managed by controlling how much the weights diverge from those of the base model. Preference optimization methods, such as Diffusion-DPO, often use a regularization factor (Beta) to control the degree of deviation. For inference, this trade-off is determined by how much of the positive weight offset η is incorporated into the negative weight α . We use the DPO algorithm to train NPO and systematically test this trade-off. Fig. 9 and Fig. 10 show examples of generated images with different parameter settings and the corresponding changes in quantitative metrics. The results indicate that choosing suitable parameters can significantly improve performance.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

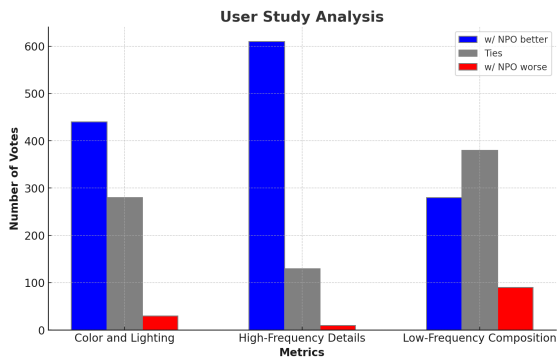


Figure 8: User study analysis.

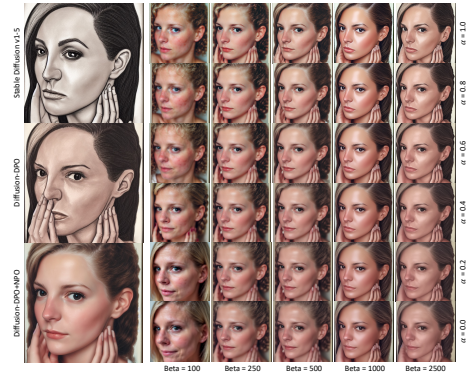


Figure 9: Visual example ablation study on hyper-parameter choice.

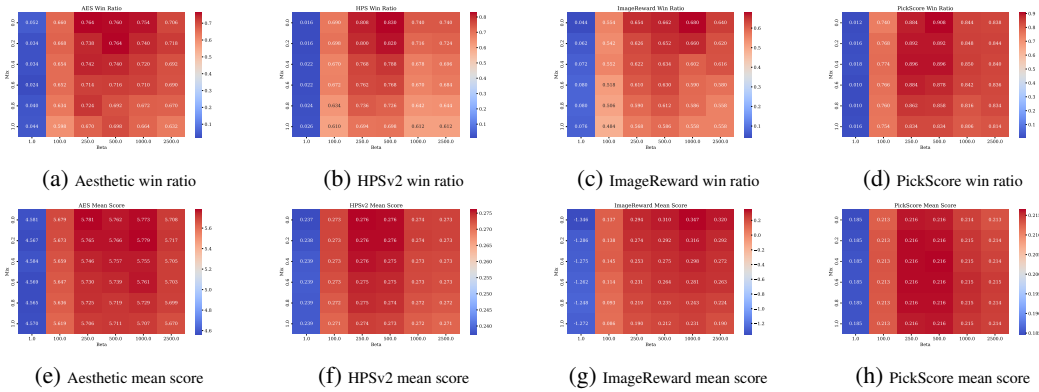


Figure 10: Heat map-based ablation study on hyper-parameter choice.

Plug-and-play. Our method is not only applicable to the original stable diffusion-based models and their fine-tuned versions optimized through preference optimization but also directly extends to high-quality stylized models fine-tuned on proprietary data. To demonstrate the versatility of our approach, we use the validation_unique dataset as our test benchmark prompts. As shown in Table 1, we observe significant improvements across various metrics. By fine-tuning the inference parameters, we enhance the performance of the DreamShaper model with 0.9B parameters, enabling it to surpass the best-performing methods on the 3B SDXL model in terms of aesthetic scores. Fig. 4 presents several comparative results, with notable improvements in structural integrity, contrast, and texture details.

5 CONCLUSIONS

In this paper, we investigate that previous preference optimization methods for diffusion models have overlooked the crucial role of unconditional/negative-conditional outputs in classifier-free guidance. We innovatively propose the task of Negative Preference Optimization as a plug-and-play inference enhancement strategy to achieve better preference-aligned generation. We summarize existing preference optimization training strategies and provide a straightforward but effective adaptation for Negative Preference Optimization. Extensive experimental results validate the effectiveness of Negative Preference Optimization.

Limitations: Diffusion-NPO requires the storage and loading of two different weight offsets for inference, which results in a higher storage cost. However, fortunately, preference optimization can typically be trained with LoRA, which requires only a minimal amount of additional storage.

REFERENCES

- 540
541
542 Donghoon Ahn, Hyoungwon Cho, Jaewon Min, Wooseok Jang, Jungwoo Kim, SeonHwa Kim,
543 Hyun Hee Park, Kyong Hwan Jin, and Seungryong Kim. Self-rectifying diffusion sampling with
544 perturbed-attention guidance. *arXiv preprint arXiv:2403.17377*, 2024.
- 545 Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion
546 models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- 547 Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler,
548 and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion mod-
549 els. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
550 pp. 22563–22575, 2023.
- 551 Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying
552 Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models. In
553 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
554 7310–7320, 2024.
- 555 Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models
556 on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.
- 557 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
558 *in neural information processing systems*, 34:8780–8794, 2021.
- 559 Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,
560 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Reinforcement learning for fine-
561 tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36,
562 2024.
- 563 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*
564 *arXiv:2207.12598*, 2022.
- 565 Tero Karras, Miika Aittala, Tuomas Kynkäänniemi, Jaakko Lehtinen, Timo Aila, and Samuli Laine.
566 Guiding a diffusion model with a bad version of itself. *arXiv preprint arXiv:2406.02507*, 2024.
- 567 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Ad-*
568 *vances in neural information processing systems*, 34:21696–21707, 2021.
- 569 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-
570 a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural*
571 *Information Processing Systems*, 36:36652–36663, 2023.
- 572 Youwei Liang, Junfeng He, Gang Li, Peizhao Li, Arseniy Klimovskiy, Nicholas Carolan, Jiao Sun,
573 Jordi Pont-Tuset, Sarah Young, Feng Yang, et al. Rich human feedback for text-to-image genera-
574 tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
575 pp. 19401–19411, 2024a.
- 576 Zhanhao Liang, Yuhui Yuan, Shuyang Gu, Bohan Chen, Tiankai Hang, Ji Li, and Liang Zheng.
577 Step-aware preference optimization: Aligning preference with denoising performance at each
578 step. *arXiv preprint arXiv:2406.04314*, 2024b.
- 579 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
580 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-
581 low instructions with human feedback. *Advances in neural information processing systems*, 35:
582 27730–27744, 2022.
- 583 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
584 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
585 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 586 Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-
587 image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023.
- 588
589
590
591
592
593

- 594 Mihir Prabhudesai, Russell Mendonca, Zheyang Qin, Katerina Fragkiadaki, and Deepak Pathak.
595 Video diffusion alignment via reward gradients. *arXiv preprint arXiv:2407.08737*, 2024.
- 596
- 597 Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
598 Wiley & Sons, 2014.
- 599 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
600 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
601 *in Neural Information Processing Systems*, 36, 2024.
- 602 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
603 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
604 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 605 Christoph Schuhmann. Laion-aesthetics. [https://laion.ai/blog/](https://laion.ai/blog/laion-aesthetics/)
606 [laion-aesthetics/](https://laion.ai/blog/laion-aesthetics/), 2022. Accessed: 2023-11-10.
- 607
- 608 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
609 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 610
- 611 Dazhong Shen, Guanglu Song, Zeyue Xue, Fu-Yun Wang, and Yu Liu. Rethinking the spatial
612 inconsistency in classifier-free diffusion guidance. In *Proceedings of the IEEE/CVF Conference*
613 *on Computer Vision and Pattern Recognition (CVPR)*, pp. 9370–9379, June 2024.
- 614 Xiaoyu Shi, Zhaoyang Huang, Fu-Yun Wang, Weikang Bian, Dasong Li, Yi Zhang, Manyuan Zhang,
615 Ka Chun Cheung, Simon See, Hongwei Qin, et al. Motion-i2v: Consistent and controllable
616 image-to-video generation with explicit motion modeling. *arXiv e-prints*, pp. arXiv–2401, 2024.
- 617 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry
618 Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video
619 data. *arXiv preprint arXiv:2209.14792*, 2022.
- 620
- 621 Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun
622 Zhou, Zipeng Qin, Yi Wang, et al. Journeydb: A benchmark for generative image understanding.
623 *Advances in Neural Information Processing Systems*, 36, 2024.
- 624 Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- 625
- 626 Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,
627 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using
628 direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
629 *and Pattern Recognition*, pp. 8228–8238, 2024.
- 630 Fu-Yun Wang, Zhaoyang Huang, Xiaoyu Shi, Weikang Bian, Guanglu Song, Yu Liu, and Hongsheng
631 Li. Animate1cm: Accelerating the animation of personalized diffusion models and adapters with
632 decoupled consistency learning. *arXiv preprint arXiv:2402.00769*, 2024.
- 633 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.
634 Human preference score v2: A solid benchmark for evaluating human preferences of text-to-
635 image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.
- 636
- 637 Xiaoshi Wu, Yiming Hao, Manyuan Zhang, Keqiang Sun, Zhaoyang Huang, Guanglu Song, Yu Liu,
638 and Hongsheng Li. Deep reward supervisions for tuning text-to-image diffusion models. *arXiv*
639 *preprint arXiv:2405.00760*, 2024.
- 640 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
641 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.
642 *Advances in Neural Information Processing Systems*, 36, 2024.
- 643 Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiabin Chen, Weihao Shen, Xiaolong Zhu, and Xiu Li.
644 Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings*
645 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8941–8951, 2024.
- 646
- 647 Yinan Zhang, Eric Tzeng, Yilun Du, and Dmitry Kislyuk. Large-scale reinforcement learning for
diffusion models. *arXiv preprint arXiv:2401.12244*, 2024.

APPENDIX

I Related Works 1

II More results 1

I RELATED WORKS

In this section, we give a brief introduction to previous efforts for diffusion-based preference optimization.

Preference datasets and reward models. Previous works including Pick-a-pic (Kirstain et al., 2023), ImageReward (Xu et al., 2024), HPSv2 (Wu et al., 2023) collect image pairs generated by diffusion models with the same prompts and label the human preference for each pair. Laion-Aesthetic (Schuhmann, 2022) asks people to rate their preference for real images from 1 to 10. They then train the preference score models based on the preference label collected. These works lay a solid foundation for future human preference optimization works in diffusion models.

Differentiable reward. Some works including DRaFT (Clark et al., 2023), AlignProp (Prabhudesai et al., 2023), and ReFL (Xu et al., 2024) directly feed the generated images into pre-trained ImageReward models and update the generative model through the gradient of differentiable reward model. These works are straightforward and effective. However, due to the imperfection of reward models, these methods typically have reward leakage. For example, they may generate over-saturated images to cheat higher scores.

Reinforcement learning. Some works including DDPO (Black et al., 2023), and DPOK (Fan et al., 2024) propose to perceive the diffusion denoising process as a Markov decision process and apply the reinforcement learning algorithms for preference alignment. Some works (Zhang et al., 2024) scale up the training for better performance. Generally, they apply PPO or the variants for training.

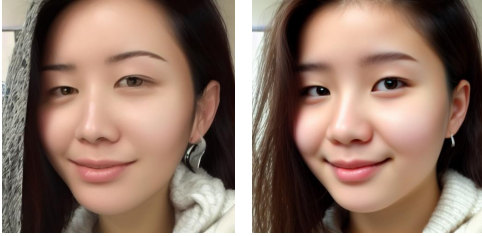
Direct Preference Optimization. Diffusion-DPO (Wallace et al., 2024) proposes a simulation-free training objective that enables direct preference optimization on preference-labeled image pairs. D3PO (Yang et al., 2024) and SPO (Liang et al., 2024b) combine reinforcement learning and direct preference optimization without the requirement to know specific score values for training.

II MORE RESULTS

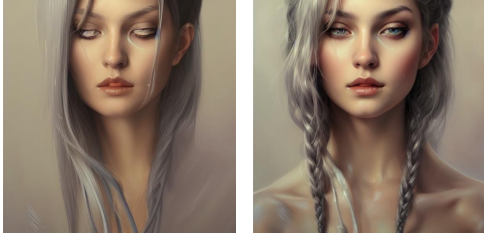
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

Stable Diffusion

Prompt: "A beautiful 25 year old whos mother is from hong kong and father from turkey"



Prompt: "a woman in a silver suit with a ponytail, a detailed painting by WLOP, trending on Artstation, fantasy art, detailed painting, artstation hd, high detail"



Prompt: "A house in the style of Escher"



Prompt: "Watercolour painting of an orange cat"



Prompt: "Milim, pink hair, that awesome time i got reincarnated as a slime"



Prompt: "Hyperrealistic full length portrait of gorgeous goddess standing in field full of flowers ... (over 30 words)"



Prompt: "female face, blue jet green eyes, long hair, slant eyes, cheeky cheeks, smiling, carefree, ... (over 20 words)"



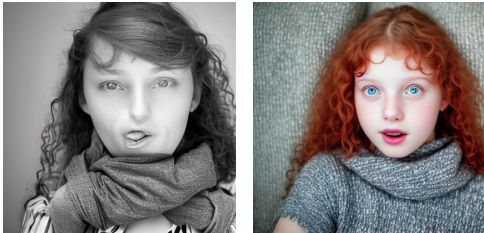
Prompt: "A giant eagle monster art"



Prompt: "An anime woman"



Prompt: "Preteen girls with no underwear neither other clothes in a sofa with a childish faces ... (over 30 words)"



w/o NPO

w/ NPO

w/o NPO

w/ NPO

Figure 11: Comparison on Stable Diffusion 1.5.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

DreamShaper

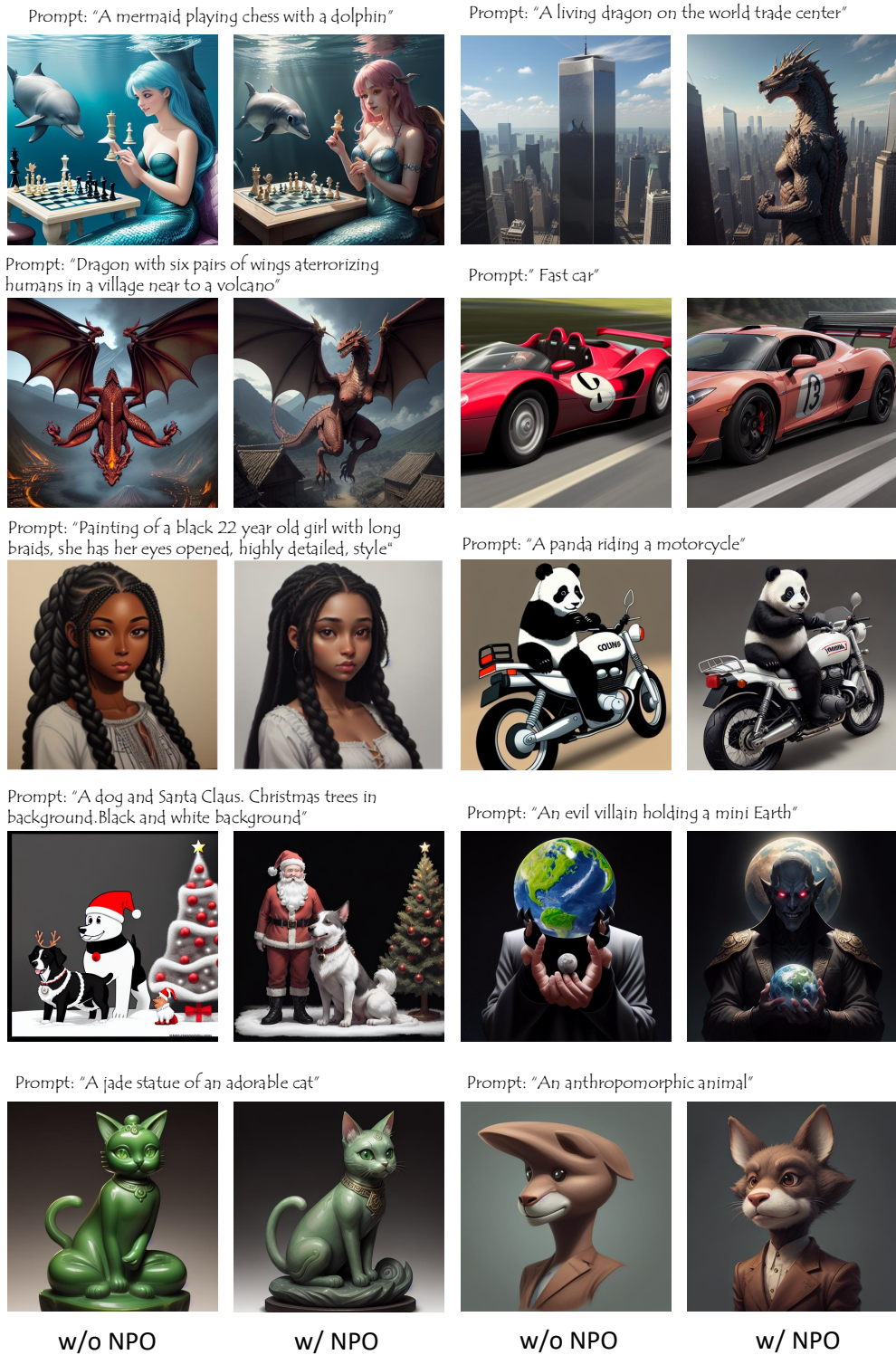
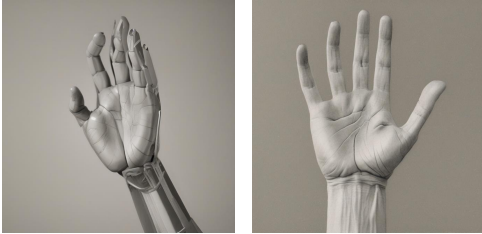


Figure 12: Comparison on DreamShaper.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

Stable Diffusion XL

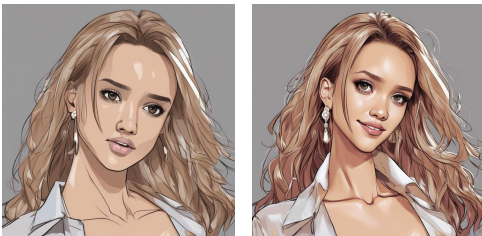
Prompt: "Human palm"



Prompt: "A cat, fat, chubby, very fine wispy and extremely long swirly wavy fur ... (over 30 words)"



Prompt: "Jessica alba, anime style"



Prompt: "LeBron James slam dunking the planet saturn through its own rings"



Prompt: "A woman with blue eyes"



Prompt: "A gijinka black cat sushi chef"



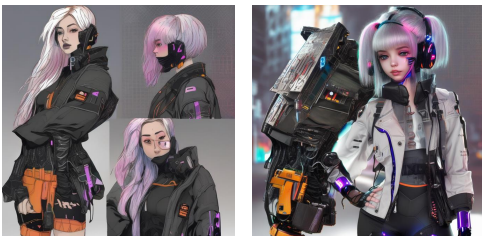
Prompt: "A boss screaming at his employee for not working on the weekend by vincent van gogh"



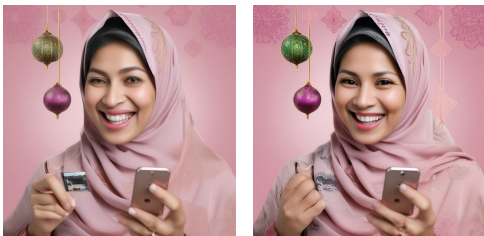
Prompt: "Concept art, Disney, really crazy creature, colored pencils, cute, very creative drawing. ... (over 30 words)"



Prompt: "A 20 yo girl in cyberpunk outfit"



Prompt: "Realistic photo with a light pink background color in various shades, a middle-aged ... (over 30 words)"



w/o NPO

w/ NPO

w/o NPO

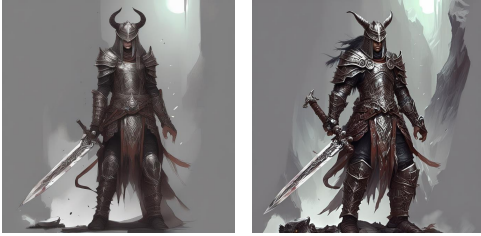
w/ NPO

Figure 13: Comparison on Stable Diffusion XL.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Diffusion-DPO

Prompt: "Fantasy warrior"



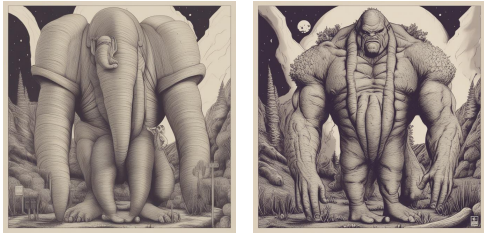
Prompt: "Wild man with a bronze axe, ring armor and furs, wielding a shield"



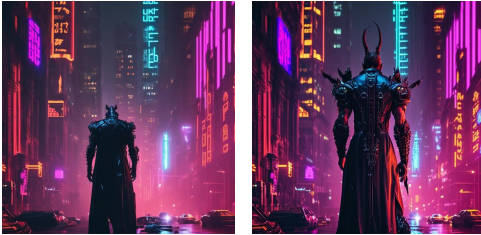
Prompt: "Realistic photo of 8 year old girl chino kafuu from is the order a rabbit, cosplay, full body"



Prompt: "Big gorilla"



Prompt: "God Hades in Gotham like city, cyberpunk, up close, cinematic, neon"



Prompt: "Sunset reflecting on a crystal ball, factory filled with android girls"



Prompt: "Smooth shading"



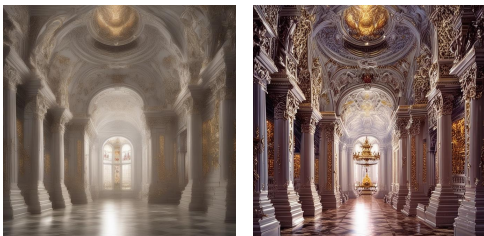
Prompt: "An attractive young woman rolling her eyes"



Prompt: "A blue car"



Prompt: "Heaven"



w/o NPO

w/ NPO

w/o NPO

w/ NPO

Figure 14: Comparison on Diffusion-DPO.

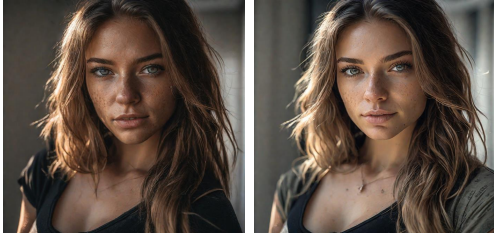
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Diffusion-SPO

Prompt: "Boy climbing into oven, lewd"



Prompt: "Rachel Amber:1.5 wearing a black skirt. Thin body type, Young face, Sony Alpha A7 III, ... (over 30 words)"



Prompt: "Japanese children ballet school"



Prompt: "Photorealistic style, photorealistic pope francis wearing drip footwear, drip tenis"



Prompt: "Michael jordan against bruce lee The straight blast round kick in the air nba basketball ball ... (over 30 words)"



Prompt: "Random girl hugs Henry Cavill superman"



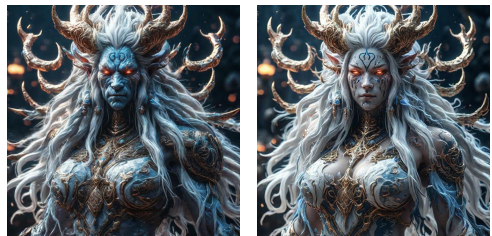
Prompt: "Highly detailed realistic photograph of a hand"



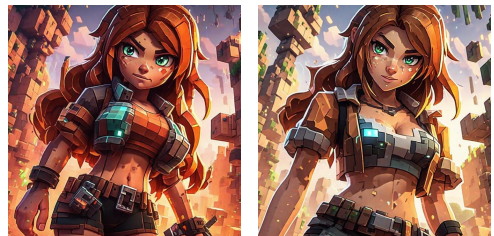
Prompt: "Portrait of a 26yr white woman, hyper-detailed, extremely ashamed, soft skin"



Prompt: "3d render of an ultrarealistic creature design, ONI entity with white long flowing hair"



Prompt: "A hot female Alex from Minecraft"



w/o NPO

w/ NPO

w/o NPO

w/ NPO

Figure 15: Comparison on Diffusion-SPO.