RETHINKING PREFERENCE ALIGNMENT FOR DIFFUSION MODELS WITH CLASSIFIER-FREE GUIDANCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Aligning large-scale text-to-image diffusion models with nuanced human preferences remains a significant challenge. Direct preference optimization (DPO), while efficient and effective, often suffers from generalization gap in large-scale finetuning. We take inspiration from test-time guidance techniques and view preference alignment as a variant of classifier-free guidance (CFG), where a finetuned preference model serves as an external control signal. This perspective yields a simple and effective method that improves alignment with human preferences. To further improve generalization, we decouple preference learning into two modules trained on positive and negative samples, whose combination at inference can yield a more effective alignment signal. We quantitatively and qualitatively validate our approach on Stable Diffusion 1.5 and Stable Diffusion XL using standard image preference datasets such as Pick-a-Pic v2 and HPDv3.

1 Introduction

Diffusion models (Ho et al., 2020; Song & Ermon, 2019; Song et al., 2021) are one of the most prevalent generative models for high-fidelity text-to-image (T2I) synthesis (Podell et al., 2023; Saharia et al., 2022). These models are typically trained from Internet-scale datasets which, due to the tremendous scale, are not carefully curated. A diffusion model pretrained on these datasets therefore deviate from what humans (in the majority voting sense) truly prefer in aspects such as aesthetic and instruction following (Kirstain et al., 2023).

The same issue is well studied in the field of large language model (LLM), in which naïvely pretrained LLMs without any post-training steps do not follow instructions and are not able to chat naturally with human (Ouyang et al., 2022). Typical approaches to align LLMs with human preference for LLMs include 1) reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022), which demands a reward model pretrained on a preference dataset and requires careful hyperparameter tuning, and 2) direct preference optimization (DPO) (Rafailov et al., 2023b), the simpler alternative that bypasses reward modeling by essentially treating the alignment problem as a binary classification problem on positive-negative preference pairs. This simple solution of DPO can be easily adapted for aligning diffusion models with human preference (Wallace et al., 2024) and has been widely used in applications other than T2I synthesis (Wang et al., 2023; Blattmann et al., 2023; Wu et al., 2023a; Khachatryan et al., 2023). Nevertheless, DPO is generally considered less robust compared to RLHF: it is prone to overfitting, may produce non-smooth predictions on out-of-distribution text prompts, and even exhibit catastrophic forgetting behaviors (Lin et al., 2024). While one may include in either the whole pretraining dataset or just the prompt set to regularize models, access to these pretraining sets is typically infeasible for large-scale models.

Since direct preference with supervised learning is often prone to overfitting issues, we instead take inspiration from inference-time techniques for diffusion model adaptation. Specifically, we observe that classifier free guidance (CFG) (Dhariwal & Nichol, 2021), the standard approach for sampling from conditional diffusion models by linearly combining between unconditional and conditional predictions, can be viewed as tempering the potentially overfitted conditional model with the more generalizable unconditional prior. Since the posterior distribution obtained through CFG typically exhibits strong performance, and the alignment objective from the control as inference perspective (Levine, 2018) is likewise to obtain a posterior distribution (Rafailov et al., 2023a), we are led to ask: *can CFG be adopted to address the diffusion alignment problem?*

Motivated by this question, we view a finetuned diffusion model as the diffusion model conditioned on a virtual control signal from the preference dataset, while the base model serves as the unconditional model. From this perspective, sampling from the aligned diffusion model naturally becomes a CFG-style inference process, which gives rise to our first method, Preference-Guided Diffusion (PGD). Adopting this CFG perspective further suggests that finetuning should resemble conditional diffusion training, which does not rely on positive–negative pairs but instead uses the standard diffusion loss. To implement this idea, we finetune two models independently, one that generates positive samples and another that generates negative samples, and combine them at inference through CFG-style composition. We refer to this variant as contrastive Preference-Guided Diffusion (cPGD). In experiments, PGD and cPGD consistently outperform vanilla Diffusion-DPO. Notably, both methods achieve Pareto improvements, simultaneously yielding higher reward, lower FID, and greater diversity in the generated samples. Moreover, the approach in principle produces a transferable plug-and-play module that, once trained on a base diffusion model, can be reused to align others.

In summary, our contributions are

- We propose to alleviate the generalization issue in Diffusion-DPO by treating diffusion model alignment as a special case of CFG-style inference.
- We introduce Preference-Guided Diffusion (PGD), which aligns the generated distribution with human preference through CFG-style guidance at inference time.
- We extend this view by considering finetuning as conditional diffusion training and propose contrastive PGD (cPGD).
- We empirically demonstrate that both variants achieve Pareto improvements over the Diffusion-DPO baseline.

2 RELATED WORK

Alignment with human preference. Preference optimization has become central to aligning large generative models with human expectations. In large language models, reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022) is the dominant framework, relying on a reward model trained from pairwise human preferences (Christiano et al., 2017; Stiennon et al., 2020). While effective, RLHF requires careful hyperparameter tuning in both the reward model and reinforcement learning stages. In contrast, direct preference optimization (DPO) (Rafailov et al., 2023b), its diffusion-specific extension Diffusion-DPO, and several related alternatives (Azizzadenesheli et al., 2023; Xu et al., 2024a; Lin et al., 2024) offer a simpler approach: directly finetuning the model with a logistic regression objective on preference pairs, thereby eliminating the need for an explicit reward model. However, DPO methods are often less competitive than RLHF (Ouyang et al., 2022), a limitation also observed in recent adaptations of preference optimization to text-to-image diffusion models (Black et al., 2023; Lee et al., 2023; Black et al., 2024; Fan et al., 2023; Xu et al., 2024b; Clark et al., 2024; Prabhudesai et al., 2023; Wallace et al., 2024; Li et al., 2024; Yang et al., 2024; Zhu et al., 2025). Building on this line of work, we propose a Diffusion-DPO variant that reformulates preference alignment as inference-time guidance to improve generalization.

Guidance in diffusion models. Controlling diffusion models can be broadly categorized into *fine-tuning approaches* and *inference-time guidance approaches*. Fine-tuning methods adapt model parameters to inject conditioning signals or domain knowledge. Representative examples include DreamBooth (Ruiz et al., 2023), which personalizes text-to-image models with subject-specific data, and other adapter- or LoRA-style techniques (Hu et al., 2021; Gal et al., 2022). While effective, such methods require additional training and may incur overfitting or catastrophic forgetting when data is limited. In contrast, inference-time guidance requires no additional training and modifies the sampling process to incorporate conditioning. Classifier guidance (Dhariwal & Nichol, 2021) uses the gradient of an external classifier, but can lead to distributional shifts. Classifier-free guidance (CFG) (Ho & Salimans, 2022) avoids this by training with randomly dropped conditions and linearly combining between unconditional and conditional predictions at inference, and has since become the de facto standard for controllable text-to-image generation. Numerous extensions build on this principle, e.g., language-model-based steering (Nichol et al., 2021), attention-based semantic guidance (Chefer et al., 2023), or plug-and-play conditioning modules (Liu et al., 2023). Our work draws direct inspiration from inference-time guidance. Instead of conditioning on textual prompts

or class labels, we extend the CFG principle to *preference alignment*, treating human preference as a conditioning signal that can be injected at inference to steer generation toward preferred outputs.

3 PRELIMINARIES

3.1 DIFFUSION MODELS

Diffusion models are a category of generative models that generate samples by sequentially denoising noisy samples. Specifically, a diffusion model defines a noising (forward) process

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{\alpha_t} \, \mathbf{x}_{t-1}, \, (1 - \alpha_t) \mathbf{I}). \tag{1}$$

where $\{\beta_t\}_{t=1}^T$ and $\alpha_t=1-\beta_t$, $\bar{\alpha}_t=\prod_{s=1}^t \alpha_s$ are the noise schedule, typically set such that $q(x_T|x_0)\approx \mathcal{N}(0,I)$. Sampling from this diffusion model is through the denoising (reverse) process:

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_{t}) = \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{x}_{t}, t), \sigma_{t}^{2}\mathbf{I}), \qquad \boldsymbol{\mu}_{\theta}(\mathbf{x}_{t}, t) = \frac{1}{\sqrt{\alpha_{t}}} \left(\mathbf{x}_{t} - \frac{\beta_{t}}{\sqrt{1 - \bar{\alpha}_{t}}} \epsilon_{\theta}(\mathbf{x}_{t}, t)\right), \quad (2)$$

with σ_t^2 set to the posterior variance $\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$.

Training a diffusion model amounts to simply minimize the diffusion loss

$$\mathcal{L}_{\text{DDPM}}(\theta) = \mathbb{E}_{t,\mathbf{x}_0,\epsilon} \left[w(t) \left\| \epsilon - \epsilon_{\theta} \left(\sqrt{\bar{\alpha}_t} \, \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \, \epsilon, \, t \right) \right\|_2^2 \right]$$
 (3)

where w(t) is a weighting scalar with one of the common choices being w(t)=1. Such as objective is equivalent to matching the model output $\epsilon_{\theta}(x,t)$ with the ground-truth score function $\nabla \log p_t(x)$ and we therefore use $\nabla \log \pi(x,t;\theta)$ interchangeably with $\epsilon_{\theta}(x,t)$.

3.2 DIRECT PREFERENCE OPTIMIZATION

Given a preference dataset $\mathcal{D} = \{(x_+^{(i)}, x_-^{(i)}, c)\}_{i=1}^N$ where x_+ and x_- are positive and negative samples (respectively) and c is the text prompt, direct preference optimization (DPO) (Rafailov et al., 2023a) aims to perform logistic regression with the relative log-odds:

$$L_{\text{DPO}} = - \underset{(x_+, x_-, c) \sim \mathcal{D}}{\mathbb{E}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(x_+|c)}{\pi_{\text{ref}}(x_+|c)} - \beta \log \frac{\pi_{\theta}(x_-|c)}{\pi_{\text{ref}}(x_-|c)} \right) \right]$$
(4)

where $\sigma(\cdot)$ is the sigmoid function. Such an objective assumes the implicit reward model with the probabilities π_{θ} and π_{ref} (with normalization constant Z): $r(x,c) = \beta \log \frac{\pi_{\theta}(x|c)}{\pi_{\text{ref}}(x|c)} + \log Z$ and the Bradley-Terry (BT) preference model (Bradley & Terry, 1952) given the optimal policy p^* :

$$p(x_{+} \succ x_{-}|c) = \sigma \left(\beta \log \frac{\pi^{*}(x_{+}|c)}{\pi_{\text{ref}}(x_{+}|c)} - \beta \log \frac{\pi^{*}(x_{-}|c)}{\pi_{\text{ref}}(x_{-}|c)}\right).$$
 (5)

In diffusion models, x_+ and x_- correspond to entire sample trajectories from t=T to t=0 if DPO is applied in the most naïve way. This is computationally intractable, since evaluating the likelihood ratio of whole trajectories requires integrating over all intermediate noise steps. Diffusion-DPO (Wallace et al., 2024) alleviates this by applying Jensen's inequality to derive an upper bound on the trajectory-level DPO loss. Specifically, the joint log-likelihood of a trajectory is decomposed into a sum of per-step transition log-likelihoods, leading to a simple transition-wise training objective:

$$L = - \underset{\substack{(x_0^+, x_0^-, c) \sim \mathcal{D} \\ t \sim U[1, T]}}{\mathbb{E}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(x_+^{(t-1)} \mid x_+^{(t)}, c)}{\pi_{\text{ref}}(x_+^{(t-1)} \mid x_+^{(t)}, c)} - \beta \log \frac{\pi_{\theta}(x_-^{(t-1)} \mid x_-^{(t)}, c)}{\pi_{\text{ref}}(x_-^{(t-1)} \mid x_-^{(t)}, c)} \right) \right], \quad (6)$$

where $\pi_{\theta}(x_{t-1} \mid x_t, c)$ denotes the one-step reverse diffusion transition probability under model π_{θ} . With the log-likelihood approximated by diffusion losses, we obtain the final form:

$$L_{\text{Diff-DPO}}(\theta) = - \underset{\substack{(x_0^+, x_0^-, c) \sim \mathcal{D} \\ t \sim U[1, T]}}{\mathbb{E}} \left[\log \sigma \left(-\beta T \, \omega(\lambda_t) \left(\| \epsilon^+ - \epsilon_{\theta}(x_+^{(t)}, t, c) \|_2^2 - \| \epsilon^+ - \epsilon_{\text{ref}}(x_+^{(t)}, t, c) \|_2^2 \right) \right]$$

$$-\|\epsilon^{-} - \epsilon_{\theta}(x_{-}^{(t)}, t, c)\|_{2}^{2} + \|\epsilon^{-} - \epsilon_{\text{ref}}(x_{-}^{(t)}, t, c)\|_{2}^{2})\right].$$
 (7)



Figure 1: Comparison of base, DPO, and PGD: PGD retains base fidelity while leveraging DPO-learned preferences.

Figure 2: Illustration of cPGD. pSFT and nSFT denote inference with the model finetuned on positive and negative samples, respectively.

3.3 CONDITIONAL GENERATION AND CLASSIFIER-FREE GUIDANCE

With a conditional diffusion model trained on the dataset $\{(x_i, c_i)\}_{i=1}^N$, the inference process typically adopts classifier-free guidance (CFG) (Ho & Salimans, 2022) that samples instead with the composed score estimate:

$$\hat{\epsilon}(\mathbf{x}_t, t, \mathbf{c}) = \epsilon_u + w \cdot (\epsilon_c - \epsilon_u), \tag{8}$$

where $\epsilon_u = \epsilon_\theta(\mathbf{x}_t, t, \varnothing), \epsilon_c = \epsilon_\theta(\mathbf{x}_t, t, \mathbf{c})$ are the unconditional and conditional score estimate (respectively), w is a positive guidance weight that is usually greater than 1, and \varnothing is the null condition. In practice, ϵ_u is trained by setting the embedding of the condition input to zero. The CFG inference process approximately generates samples from the posterier distribution $p(x)p^w(c|x)$ with p(x) being the prior distribution, or equivalently in log-likelihood, $\log p(x) + w \log p(c|x)$.

4 Method

4.1 Preference-Guided Inference

Let $\pi_{\rm ref}$ be a reference policy and $\pi_{\rm DPO}$ be a DPO-tuned policy. By treating the DPO-tuned policy as $\pi(x|\mathcal{D})$ and the reference policy as $\pi(x|\mathcal{D})$ as in CFG, we immediately obtain the CFG-style score function for inference, for which we term preference-guided diffusion (PGD):

$$\nabla \log \pi_{PGD}(x) = \nabla \log \pi_{ref}(x) + w \Big(\nabla \log \pi_{DPO}(x) - \nabla \log \pi_{ref}(x) \Big), \tag{9}$$

where the guidance weight w determines the trade-off between our confidence on the reward and other metrics such as prior preservation and sample diversity. Since $\pi_{\rm ref}$ can be understood as some prior pretrained on unlabeled datasets, once we have trained $\nabla \log \pi_{\rm DPO}(x)$, we are able to virtually align any other base model $\pi'_{\rm ref}(x)$ by simply replacing $\pi_{\rm ref}(x)$ with it.

4.2 Contrastive PGD as Dynamically-Reweighted Guidance

The connection between CFG and diffusion model alignment prompts us to think whether finetuning should also be done in a way similar to conditional diffusion model training, which directly encourage the negative score functions to point towards the data points. However, our preference dataset contains both positive samples and negative ones. These negative samples act as "repelling" forces that pushes the negative score function away from them. Inspired by that much of alignment can be turned into an inference-time manner, we propose to postpone this "repelling" behavior to inference-time as well. Specifically, we finetune another copy of the base model so that it generates negative samples. Formally speaking, with \mathcal{D}_+ representing the set of positive samples and

 \mathcal{D}_{-} the set of negative ones, we independently finetune two models (with parameters θ_{+} and θ_{-} , respectively) with diffusion losses:

$$L_{\text{pos}}(\theta_{+}) = \mathbb{E}_{\epsilon \sim \mathcal{N}(0,I), t \sim \text{Uniform}(0,1), x_{0} \sim \mathcal{D}_{+}} \left\| \epsilon - \epsilon \left(\sqrt{\bar{\alpha}_{t}} \, x_{0} + \sqrt{1 - \bar{\alpha}_{t}}, \epsilon, t; \theta_{+} \right) \right\|^{2}$$
(10)

$$L_{\text{neg}}(\theta_{-}) = \mathbb{E}_{\epsilon \sim \mathcal{N}(0,I), t \sim \text{Uniform}(0,1), x_{0} \sim \mathcal{D}_{-}} \left\| \epsilon - \epsilon \left(\sqrt{\bar{\alpha}_{t}} x_{0} + \sqrt{1 - \bar{\alpha}_{t}}, \epsilon, t; \theta_{-} \right) \right\|^{2}$$
(11)

Intuitively, the difference between two models characterizes the implicit reward model. Therefore we may write the residual parameterization $\nabla \log \pi_{\text{finetuned}}(x,t) = \nabla \log \pi(x,t;\theta_+) - \nabla \log \pi(x,t;\theta_-) + \nabla \log \pi_{\text{ref}}(x,t)$. It follows that the resulting PGD formulation, to which we refer with contrastive PGD (cPGD), is

$$\nabla \log \pi_{\text{PGD}}(x,t) = \nabla \log \pi_{\text{ref}}(x,t) + w \Big(\nabla \log \pi(x,t;\theta_{+}) - \nabla \log \pi(x,t;\theta_{-}) \Big). \tag{12}$$

Alternative perspective of cPGD. While it may seem a bit arbitrary to replace the DPO loss on the preference dataset with two diffusion losses on positive-only and negative-only datasets respectively, cPGD essentially performs dynamic reweighting of DPO loss gradients. For simplicity, let's consider the general DPO case (without Diffusion-DPO approximations). If we plug the residual parametrization of the finetuned model into the DPO loss gradient, we observe (with $\theta = (\theta_+, \theta_-)$):

$$\nabla_{\theta} L_{\text{DPO}} = - \underset{(x_{+}, x_{-}) \sim \mathcal{D}}{\mathbb{E}} \left[\beta \sigma \left(\log \pi(x; \theta_{-}) - \log \pi(x; \theta_{+}) \right) \cdot \left(\nabla_{\theta_{+}} \log \pi(x; \theta_{+}) - \nabla_{\theta_{-}} \log \pi(x; \theta_{-}) \right) \right]. \tag{13}$$

Suppose that, for each sample pair (x_+, x_-) , we dynamically reweight the loss function by $1/\left[\beta\sigma\left(\log\pi(x;\theta_-)-\log\pi(x;\theta_+)\right)\right]$. The resulting dynamically-reweighted loss is

$$\nabla_{\theta} L_{\text{reweight}} = - \underset{(x_{+}, x_{-}) \sim \mathcal{D}}{\mathbb{E}} \left[\nabla_{\theta_{+}} \log \pi(x; \theta_{+}) - \nabla_{\theta_{-}} \log \pi(x; \theta_{-}) \right]$$

$$= \underset{x_{-} \sim \mathcal{D}_{-}}{\mathbb{E}} \left[\nabla_{\theta_{-}} \log \pi(x; \theta_{-}) \right] - \underset{x_{+} \sim \mathcal{D}_{+}}{\mathbb{E}} \left[\nabla_{\theta_{+}} \log \pi(x; \theta_{+}) \right]$$
(14)

which is exactly the gradient of the cPGD training objectives once we take into consideration the fact that $\log \pi$ is parameterized by a diffusion model. Wu et al. (2025) show that such reweighting can be seen as an interpolation between supervised finetuning gradients and vanilla policy gradients, and that it can be helpful to alleviate the overfitting issue due to the small scale of finetuning datasets.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Training datasets. We consider consider two datasets: 1) Pick-a-Pic v2 (Kirstain et al., 2023), which consists of approximately 900,000 image preference pairs derived from 58,000 unique prompts, and 2) HPDv3 (Ma et al., 2025), which comprises 1.08M text-image pairs and 1.17M annotated pairwise data. Our main experiments are done on Pick-a-Pic v2, while for ablation, we create a high-image-quality subset of HPDv3 besides the full dataset of HPDv3.

Test prompts. We consider the following prompt datasets for testing: the test split of Pick-a-Pic v2 (424 prompts), the HPDv2 test set (Wu et al., 2023b) (400 prompts), and the Parti-Prompts benchmark (Yu et al., 2022) (1,632 prompts).

Baselines. We benchmark our approaches against the following baselines: (i) SFT-Pref, a supervised fine-tuning baseline using only the preferred images; (ii) Diffusion-DPO (Wallace et al., 2024), an adaptation of the DPO method to diffusion models; (iii) Diffusion-KTO (Li et al., 2024), a variant that incorporates a Kullback–Leibler trade-off to Diffusion-DPO for unlocking the potential of leveraging readily available per-image binary signals; (iv) MaPO (Hong et al., 2024), which refines preference optimization with margin-based pairwise consistency; (v) Diffusion-NPO (Wang et al., 2025), which explicitly models negative preferences to strengthen classifier-free guidance. Additionally, we consider SPO (Liang et al., 2024), a hybrid method that trains auxiliary reward models

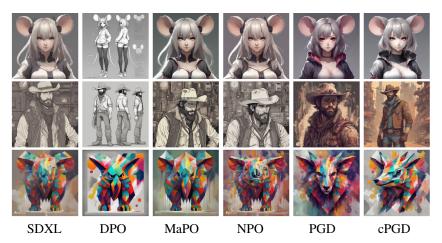


Figure 3: Comparison of preference-optimization methods on SDXL. Columns show outputs from the base model (SDXL), DPO, MaPO, NPO, PGD, and cPGD and cPGD achieves the highest rewards and is the most effective in aligning with human preference implied in the Pick-a-Pic v2 dataset.

in an online fashion during DPO finetuning; due to its online nature, we exclude SPO from the direct comparison between offline DPO variants.

Reward models. We evaluate generated images with these reward models: PickScore (PS) (Kirstain et al., 2023), HPSv2 (Wu et al., 2023b), HPSv3 (Ma et al., 2025), ImageReward (IR) (Xu et al., 2024b), CLIP Score (Radford et al., 2021), and Aesthetics Score (Aes) (Schuhmann, 2022).

Metrics. Besides the absolute reward values, we compute win rates for different methods, which is the percentage of instances where the finetuned model outperforms the base model. Since win rates are considerably more robust than absolute reward values (Wu et al., 2023b; Kirstain et al., 2023), we use win rate as our primary metric. In addition, we compute FID score and diversity score to measure the extent of prior preservation and sample diversity, respectively. Sample diversity is computed by measuring the average pairwise distances between CLIP embeddings Radford et al. (2021) of generated samples (details in Appendix A.2).

Implementation Details. We experiment with two base models: Stable Diffusion v1.5 (SD1.5) (Rombach et al., 2022) and Stable Diffusion XL base (SDXL) (Podell et al., 2023). An effective batch size of 2048 image pairs is used for all experiments. Following common practices, we set for each model a base learning rate and scale it linearly with the batch size. For SD1.5, we use AdamW optimizer with a base learning rate of 3e-8; for the larger SDXL model, we employ Adafactor optimizer with a base learning rate of 5e-9. Following Wallace et al. (2024), β is set to 3000 and 5000 for SD1.5 and SDXL, respectively. It is worth noting, however, that our effective learning rate is smaller than that of Wallace et al. (2024). In their setting, the effective learning rate is 2.048e-5, whereas ours is only 9.6e-7, an order of magnitude smaller. For PGD, we use the finetuned model with training 2000 steps and for cPGD, we use models trained with 500 steps.

5.2 RESULTS

General results. As shown in Fig. 4, Table 1 and Table 2, we find that our proposed method PGD and cPGD generally outperform the baselines in achieving higher absolute reward values and win rates for different test prompt sets and different base models. While our methods generally achieve lower Aes scores, the behavior is less indicative because our training objective is to align with the human preference implied by text-image paired datasets, while Aes is an unconditional reward model that does not take text-image alignment into consideration.

Diversity and prior preservation. We further demonstrate the tradeoffs between reward, FID (measuring prior preservation) and diversity scores in Fig. 8 and Fig. 7. The blue regions are the combinations that are strictly dominated by the performance of our methods, the boundary of which is formed by the performance resulted from different choices of guidance weights.

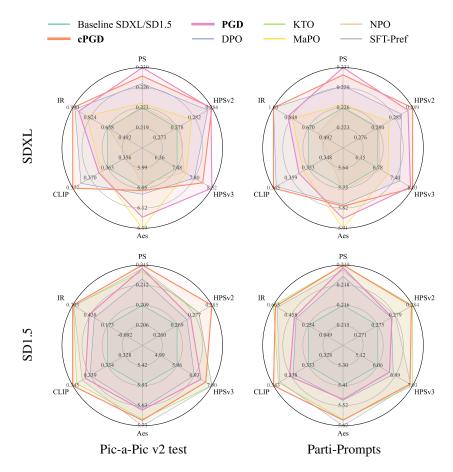


Figure 4: Overall comparison on SDXL (top) and SD1.5 (bottom). Radar axes report mean scores (higher is better): PickScore (PS), HPSv2, HPSv3, Aesthetics (Aes), CLIP, and ImageReward (IR). Polygons closer to the outer rim indicate better aggregate performance across metrics.

Table 1: Win rates of preference optimization methods against the SDXL model on the Pick-a-Pic v2 test set and the Parti-Prompts benchmark. Model checkpoints for other methods are provided by their respective authors. The 1st-best results are **bolded** and the 2nd-best results are underlined.

Base Model	Inference Strategy	PS ↑	Pick-a-Pic v2 test (424 prompts) HPSv2 ↑ HPSv3 ↑ Aes ↑ CLIP ↑ IR ↑				IR ↑	Parti-Prompts (1632 prompts) PS↑ HPSv2↑ HPSv3↑ Aes↑ CLIP↑ IR↑					
	Strategy	FS	nrsv2	пгэхэ	Aes	CLIF	IK	F5	nrsv2	nrsvs	Aes	CLIF	IK
SDXL	_	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0
	PGD	78.8	79.0	73.6	51.9	63.7	69.3	78.7	78.4	78.2	68.8	54.0	69.1
	cPGD	80.0	80.2	77.1	50.9	64.9	69.8	75.8	80.3	77.8	57.2	62.0	74.0
	NPO	58.7	59.2	69.1	52.1	37.5	53.5	55.0	56.4	62.1	55.1	39.0	51.5
DPO-SDXL	_	71.7	77.6	67.9	53.3	61.6	65.8	64.0	70.3	64.0	57.5	58.5	69.7
	PGD	83.3	85.4	85.6	59.7	62.3	73.6	80.8	83.9	81.7	67.6	57.5	76.1
DI O-SDAL	cPGD	80.9	77.6	84.7	63.9	58.7	64.9	73.9	79.2	72.8	59.4	64.1	74.4
	NPO	76.9	81.8	81.4	53.8	57.8	70.8	70.6	78.4	77.5	63.4	56.3	70.8
MaPO-SDXL	_	55.9	65.3	61.8	68.2	50.2	68.2	52.0	64.4	58.5	72.4	48.2	65.0
	PGD	80.4	81.6	81.6	75.7	51.4	72.2	78.9	77.8	79.6	77.8	53.6	72.7
	cPGD	77.4	78.8	72.9	69.1	59.4	72.4	72.5	81.1	76.9	70.1	58.2	74.4
SPO-SDXL	_	89.4	83.0	96.0	81.8	33.3	78.8	87.8	85.5	92.2	88.1	31.7	74.9
	PGD	92.2	86.1	96.5	82.1	42.5	81.4	91.4	87.7	93.9	88.4	48.0	77.3
	cPGD	92.9	88.4	96.7	78.8	53.8	84.4	92.0	90.3	93.9	83.6	50.4	81.3

PGD vs. cPGD. We find that cPGD is generally better on SD1.5 but comparable on SDXL. We hypothesize that such behavior is due to the distribution shift between the image distributions of the preference datasets and that of the base model. As the images in the preference datasets we used are generally better than those from SD1.5, the dynamic reweighting mechanism used in cPGD helps generalization in this case.

Transfer to other base models. Inspired by the plug-and-play nature of our approach, we experiment with aligning base models that are finetuned with alternative DPO variants on the same preference datasets using our PGD/cPGD-finetuned modules (*e.g.*, the second mega-row "DPO-SDXL"

Table 2: Win rates of preference optimization methods against the SD1.5 model on the Pick-a-Pic v2 test set and the Parti-Prompts benchmark. Model checkpoints for other methods are provided by their respective authors. The 1st-best results are **bolded** and the 2nd-best results are underlined.

Base Model	Inference		Pick-a	-Pic v2 test (424 prom	ipts)	Parti-Prompts (1632 prompts)				ots)		
Dasc Model	Strategy	PS ↑	HPSv2↑	HPSv3↑	Aes ↑	CLIP ↑	IR ↑	PS ↑	HPSv2↑	HPSv3 ↑	Aes ↑	CLIP ↑	IR ↑
SD1.5	-	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0	50.0
	PGD	78.3	71.2	67.9	62.3	58.5	63.7	68.0	65.0	59.6	58.1	55.0	56.9
	cPGD	76.9	71.7	71.9	63.2	59.9	72.2	66.4	76.9	68.9	68.1	58.4	71.0
DPO-SD1.5	_	76.4	67.7	66.3	65.1	55.9	60.6	67.3	64.8	64.5	62.2	53.6	61.0
	PGD	79.2	70.3	66.3	66.3	59.4	63.2	74.2	67.4	65.0	63.1	55.5	62.9 67.3
	cPGD	79.0	81.8	75.5	71.7	62.3	76.2	73.8	74.1	69.6	69.0	59.0	67.3
KTO-SD1.5	-	72.6	78.1	76.2	68.6	58.5	75.0	66.6	78.3	71.9	68.8	53.3	71.3
	PGD	81.6	83.3	80.2	70.3	61.1	77.4	72.1	80.3	72.8	72.4	55.1	73.8
	cPGD	76.7	80.2	75.9	70.5	60.4	74.3	66.2	77.1	69.5	68.4	55.9	72.2
SPO-SD1.5	_	71.2	63.0	64.9	68.2	38.7	61.1	68.6	61.2	64.2	71.9	37.7	61.6
	PGD	79.7	70.3	69.6	69.8	44.3	67.9	74.2	66.5	66.2	72.2	46.7	67.0
	cPGD	82.3	81.8	75.5	71.2	60.6	76.2	74.8	71.9	71.3	73.9	47.5	72.9

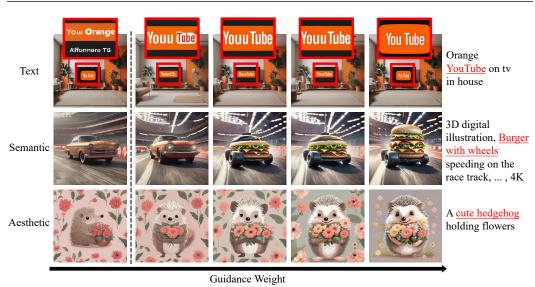
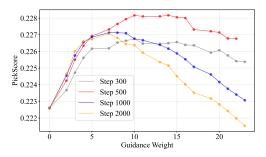


Figure 5: Qualitative effect of increasing guidance weight w (left \rightarrow right). Rows show text fidelity, semantic binding, and aesthetic style. Stronger w improves alignment and legibility up to a mid range, after which overshooting/rigidity appears.

in Table 1 demonstrate the performance when using DPO-tuned SDXL as the inference-time base model). We find that there is nearly consistent improvement compared to any original base model, which is made easy with the CFG-style inference rule in our PGD method.

Ablation on guidance weights. Qualitatively, increasing guidance weights generally yields better reward following. as shown in Fig. 5. To comprehensively quantify the effect of guidance weights in different hyperparameter settings, we measuring the performance metrics by varying the guidance weights for models finetuned on Pick-a-Pic v2 from SDXL with 300, 500, 1000 and 2000 steps. As shown in Fig. 6, we observe that increasing guidance weights from 0 to some moderate value (around 6) generally leads to better reward values for all tested models, but beyond that the model performance drops. Models finetuned with less steps exhibit less amount of performance drop, which is likely due to the regularization effect of early stopping. Furthermore, Fig. 8 and 7 shows that 1) the increase of guidance weights leads to "less natural" images and 2) if the guidance weight is beyond certain threshold, increasing the guidance weight leads to more chaotic predictions.

Dataset quality. Since the image distribution in the preference datasets can differ from the image distribution of the base models, here we investigate how methods are robust to different preference datasets, especially when the image quality differs a lot. In Table 3, we show the results of finetuning on the full HPDv3 dataset, in which the variance of image quality is great, and on a high-quality subset. We find that our methods generally performs better in both cases, but on the high-quality subset our methods unanimously outperform the baselines. In addition, in the high-quality subset case, cPGD is unanimously better than PGD. We hypothesize that this is likely due to that cPGD



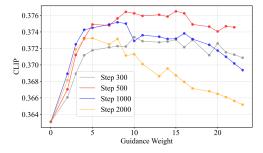


Figure 6: Effect of guidance weight w on automatic metrics (SDXL). Left: PickScore; Right: CLIP score. "Step" denotes the training steps of the guidance module. Curves rise quickly for small w,

Table 3: Impact of preference data variance on alignment performance. "Subset" refers to training on a high-quality curated subset (low variance), while "Fullset" uses the full HPDv3 dataset (high variance). The 1st-best results are in **bold** and the 2nd-best are underlined. All methods are applied to the base SDXL model.

Set	Method	PS ↑	HPSv2↑	HPSv3↑	Aes ↑	CLIP ↑	IR ↑
	SDXL	0.2226	0.2777	7.0795	6.0521	0.3631	0.6583
Subset	DPO	0.2253	0.2828	8.0327	6.1124	0.3619	0.8231
Subset	PGD	0.2257	0.2854	9.4588	6.2029	0.3637	0.8741
	cPGD	0.2276	0.2889	10.0454	6.2670	0.3646	1.0312
	SDXL	0.2226	0.2777	7.0795	6.0521	0.3631	0.6583
Fullset	DPO	0.2266	0.2847	8.2610	6.1658	0.3659	0.9179
runset	PGD	0.2285	0.2871	10.0649	6.5050	0.3644	1.0625
	cPGD	0.2273	0.2902	9.2426	6.1791	0.3734	1.1433

imposes weaker assumptions on preference pairs as θ^+ and θ^- are trained in an independent way. Despite that the high-quality subset yields more consistent observations, using the full dataset generally leads to better reward values, in part simply due to the increased number of data points.

6 DISCUSSIONS

PGD as kernel method. As shown in our experiments, PGD inference with slightly finetuned models consistently outperforms DPO methods. This behavior can be understood through the theory of neural tangent kernels (NTK) (Jacot et al., 2018). In the *lazy training regime*, i.e. when the finetuned model remains close to the reference model, we can write $\epsilon_{\rm ref} + w(\epsilon_{\rm finetuned} - \epsilon_{\rm ref}) \approx \epsilon_{\rm ref} + wK_{\rm ref}\alpha$ where $K_{\rm ref}$ is known as the NTK matrix of $\epsilon_{\rm ref}$ and α is the vector of regression coefficients, which shows that PGD is essentially kernel regression in the NTK feature space of the reference model. Because this feature space is an intrinsic and stable property of $\epsilon_{\rm ref}$, PGD inference leverages reliable features. In contrast, extended finetuning with large learning rates can push the model out of the lazy regime, causing the NTK approximation to break down and increasing the risk of overfitting on small datasets.

Inference cost. While the inference time is doubled with PGD due to the need to compute outputs with the reference model, we note that it is possible to perform distillation so that one single model learns to predict the PGD outputs, as demonstrated by many other works on diffusion model distillation (Salimans & Ho, 2022; Song et al., 2023; Meng et al., 2023). To verify this, we present our attempts in distilling a single model out of PGD in Appendix A.3.

7 Conclusion

We introduced preference-guided diffusion (PGD), a simple yet effective method that better aligns diffusion models with human preference through the lens of classifier-free guidance: the finetuned model is the guidance signal of the dataset. By further take inspiration from the training of conditional diffusion models, we propose a variant called contrastive PGD (cPGD) which parameterize the finetuned module with two models independently trained on positive and negative samples, respectively. We empirically verify the effectiveness of the proposed methods on different datasets and base models.

ETHICS STATEMENT

This work investigates preference-guided generation for text-to-image diffusion models. We did not collect new human subjects data; all experiments use publicly available datasets and prompt suites (e.g., Pick-a-Pic v2, HPDv2/3, Parti-Prompts) under their respective licenses and intendeduse policies. No personally identifiable information (PII) was processed to our knowledge. Where dataset curators provide safety filters or content flags, we follow them and do not intentionally prompt for unsafe content.

Preference signals and reward models may reflect societal biases (e.g., aesthetics tied to culture, gender, or geography). Such biases can be amplified by guidance at inference time. We therefore (i) report multiple metrics, including diversity and base-model fidelity, (ii) encourage downstream deployers to pair our method with content filters, auditing on representative user groups, and opt-out mechanisms, and (iii) commit to releasing prompts, seeds, and code sufficient for reproducibility while avoiding lists or examples that facilitate misuse (e.g., targeted impersonation or non-consensual content).

Finally, we oppose harmful uses of generative models (e.g., harassment, disinformation, infringement) and will abide by takedown requests from dataset owners within their policies.

REFERENCES

- Kamyar Azizzadenesheli, Pier Giuseppe Sessa, Kartikeya Anand, Ankesh Anand, and Maryam Fazel. Mapo: Model-agnostic preference optimization. *arXiv preprint arXiv:2310.03708*, 2023.
- Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- Sam Black, Leo Gao, Stella Biderman, Eric Hallahan, Quentin Anthony, Jason Phang, Shimao Prakash, Janelle Pfau, Shreyas Purohit, Chris Foster, et al. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- Andreas Blattmann et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Hila Chefer, Ron Mokady, Amir Bar, Amit Zohar, Roni Paiss, and Lior Wolf. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.
- Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2818–2829, 2023.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models on differentiable rewards. In *The Twelfth International Conference on Learning Representations*, 2024.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- Carles Domingo-Enrich, Michal Drozdzal, Brian Karrer, and Ricky T. Q. Chen. Adjoint matching: Fine-tuning flow and diffusion generative models with memoryless stochastic optimal control, 2025. URL https://arxiv.org/abs/2409.08861.

Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,
 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for
 fine-tuning text-to-image diffusion models. Advances in Neural Information Processing Systems,
 36:79858–79885, 2023.

- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, 2017.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In Advances in Neural Information Processing Systems (NeurIPS) Workshop on Deep Generative Models, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020.
- Jiwoo Hong, Sayak Paul, Noah Lee, Kashif Rasul, James Thorne, and Jongheon Jeong. Margin-aware preference optimization for aligning diffusion models without reference, 2024.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Levon Khachatryan et al. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023.
- Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Picka-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:36652–36663, 2023.
- Kuan-Hao Lee, Saining Xie, Han Zhang, Yiming Zhang, Chong Zhang, Luke Zettlemoyer, and Tatsunori B Hashimoto. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review. *arXiv preprint arXiv:1805.00909*, 2018.
- Shufan Li, Konstantinos Kallidromitis, Akash Gokul, Yusuke Kato, and Kazuki Kozuka. Aligning diffusion models by optimizing human utility. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Zhanhao Liang, Yuhui Yuan, Shuyang Gu, Bohan Chen, Tiankai Hang, Mingxi Cheng, Ji Li, and Liang Zheng. Aesthetic post-training diffusion models from generic preferences with step-by-step preference optimization. *arXiv preprint arXiv:2406.04314*, 2024.
- Yong Lin, Skyler Seto, Maartje Ter Hoeve, Katherine Metcalf, Barry-John Theobald, Xuan Wang, Yizhe Zhang, Chen Huang, and Tong Zhang. On the limited generalization capability of the implicit reward model induced by direct preference optimization. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 16015–16026, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.940. URL https://aclanthology.org/2024.findings-emnlp.940/.
- Zhen Liu, Tim Z. Xiao, Weiyang Liu, Yoshua Bengio, and Dinghuai Zhang. Efficient diversity-preserving diffusion alignment via gradient-informed gflownets. In *International Conference on Learning Representations*, 2025.

- Zhihao Liu, Tan Yu, Jiuxiang Gu, Ruixiang Zhang, Yi Yang, Zhangyang Wang, and Bolei Zhou.
 Plug-and-play diffusion features for text-driven image-to-image translation. *arXiv preprint arXiv:2304.02883*, 2023.
 - Yuhang Ma, Yunhao Shui, Xiaoshi Wu, Keqiang Sun, and Hongsheng Li. Hpsv3: Towards wide-spectrum human preference score, 2025. URL https://arxiv.org/abs/2508.03789.
 - Chenlin Meng, Yang Song, Jiaming Song, Jiaming Wu, Prafulla Dhariwal, Alexey Nichol, Xinyang Chen, Jascha Sohl-Dickstein, and Stefano Ermon. On distillation of guided diffusion models. *arXiv* preprint arXiv:2306.05544, 2023.
 - Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
 - Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
 - Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023.
 - Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-image diffusion models with reward backpropagation, 2023.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.
 - Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023a. URL https://arxiv.org/abs/2305.18290.
 - Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023b.
 - Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
 - Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
 - Chitwan Saharia, William Chan, Saurabh Saxena, Lala Lit, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Raphael Gontijo-Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
 - Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*, 2022.
 - Christoph Schuhmann. Laion-aesthetics. https://laion.ai/blog/laion-aesthetics/, 2022. Accessed: 2023 11- 10.

- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv* preprint arXiv:2010.02502, 2020.
 - Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
 - Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL https://openreview.net/forum?id=PxTIG12RRHS.
 - Yang Song, Chenlin Meng, and Stefano Ermon. Consistency models. In *International Conference on Machine Learning*, 2023.
 - Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
 - Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8228–8238, 2024.
 - Fu-Yun Wang, Yunhao Shui, Jingtan Piao, Keqiang Sun, and Hongsheng Li. Diffusion-npo: Negative preference optimization for better preference aligned generation of diffusion models, 2025. URL https://arxiv.org/abs/2505.11245.
 - Yifan Wang et al. Mesh-rft: Reframing text-to-image diffusion for 3d mesh generation. *arXiv* preprint arXiv:2303.XXXX, 2023.
 - Tianxing Wu, Yijun Ge, Xintao Zhang, et al. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *ICCV*, 2023a.
 - Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis, 2023b. URL https://arxiv.org/abs/2306.09341.
 - Xuehai Wu, Jinghan Zhang, Haonan Dong, Yujia Li, Wenlong Hu, Wayne Xin Zhao, and Ji-Rong Wen. On the generalization of sft: A reinforcement learning perspective with reward rectification. arXiv preprint arXiv:2508.05629, 2025. URL https://arxiv.org/abs/2508.05629.
 - Frank F Xu, Xuechen Li, Faisal Ladhak, Esin Durmus, Jason Wei, Xiang Lorraine Chen, and Tatsunori B Hashimoto. Implicit preference optimization: Aligning language models without a reward model. *arXiv preprint arXiv:2402.00856*, 2024a.
 - Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36, 2024b.
 - Shentao Yang, Tianqi Chen, and Mingyuan Zhou. A dense reward view on aligning text-to-image diffusion with preference. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 55998–56032, 2024.
 - Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for contentrich text-to-image generation. *Trans. Mach. Learn. Res.*, 2022.
 - Huaisheng Zhu, Teng Xiao, and Vasant G Honavar. DSPO: Direct score preference optimization for diffusion model alignment. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=xyfb9HHvMe.