

Diffusion Model Alignment Using Direct Preference Optimization

Bram Wallace¹ Meihua Dang² Rafael Rafailov² Linqi Zhou² Aaron Lou²
Senthil Purushwalkam¹ Stefano Ermon² Caiming Xiong¹ Shafiq Joty¹
Nikhil Naik¹

¹Salesforce AI, ²Stanford University

{b.wallace, spurushwalkam, cxiong, sjoty, nnaik}@salesforce.com

{mhdang, rafailov, lzhou907, aaronlou}@stanford.edu {ermon}@cs.stanford.edu

Abstract

Large language models (LLMs) are fine-tuned using human comparison data with Reinforcement Learning from Human Feedback (RLHF) methods to make them better aligned with users' preferences. In contrast to LLMs, human preference learning has not been widely explored in text-to-image diffusion models; the best existing approach is to fine-tune a pretrained model using carefully curated high quality images and captions to improve visual appeal and text alignment. We propose Diffusion-DPO, a method to align diffusion models to human preferences by directly optimizing on human comparison data. Diffusion-DPO is adapted from the recently developed Direct Preference Optimization (DPO) [33], a simpler alternative to RLHF which directly optimizes a policy that best satisfies human preferences under a classification objective. We re-formulate DPO to account for a diffusion model notion of likelihood, utilizing the evidence lower bound to derive a differentiable objective. Using the Pick-a-Pic dataset of 851K crowdsourced pairwise preferences, we fine-tune the base model of the state-of-the-art Stable Diffusion XL (SDXL)-1.0 model with Diffusion-DPO. Our fine-tuned base model significantly outperforms both base SDXL-1.0 and the larger SDXL-1.0 model consisting of an additional refinement model in human evaluation, improving visual appeal and prompt alignment. We also develop a variant that uses AI feedback and has comparable performance to training on human preferences, opening the door for scaling of diffusion model alignment methods.

1. Introduction

Text-to-image diffusion models have been the state-of-the-art in image generation for the past few years. They are typically trained in a single stage, using web-scale datasets of text-image pairs by applying the diffusion objective. This

stands in contrast to the state-of-the-art training methodology for Large Language Models (LLMs). The best performing LLMs [28, 48] are trained in two stages. In the first (“pretraining”) stage, they are trained on large web-scale data. In the second (“alignment”) stage, they are fine-tuned to make them better aligned with human preferences. Alignment is typically performed using supervised fine-tuning (SFT) and Reinforcement Learning from Human Feedback (RLHF) using preference data. LLMs trained with this two-stage process have set the state-of-the-art in language generation tasks and have been deployed in commercial applications such as ChatGPT and Bard.

Despite the success of the LLM alignment process, most text-to-image diffusion training pipelines do not incorporate learning from human preferences. Several models [9, 35, 36], perform two-stage training, where large-scale pretraining is followed by fine-tuning on a high-quality text-image pair dataset to strategically bias the generative process. This approach is much less powerful and flexible than the final-stage alignment methods of LLMs. Recent approaches [6, 7, 11, 31] develop more advanced ways to tailor diffusion models to human preferences, but none have demonstrated the ability to stably generalize to a fully open-vocabulary setting across an array of feedback. RL-based methods [6, 11] are highly effective for limited prompt sets, but their efficacy decreases as the vocabulary expands. Other methods [7, 31] use the pixel-level gradients from reward models on generations to tune diffusion models, but suffer from mode collapse and can only be trained on a relatively narrow set of feedback types.

We address this gap in diffusion model alignment for the first time, developing a method to directly optimize diffusion models on human preference data. We generalize Direct Preference Optimization (DPO) [33], where a generative model is trained on paired human preference data to implicitly estimate a reward model. We define a notion of data likelihood under a diffusion model in a novel formulation and derive a simple but effective loss resulting in stable



Figure 1. We develop Diffusion-DPO, a method based on Direct Preference Optimization (DPO) [33] for aligning diffusion models to human preferences by directly optimizing the model on user feedback data. After fine-tuning on the state-of-the-art SDXL-1.0 model, our method produces images with exceptionally high visual appeal and text alignment, samples above.

and efficient preference training, dubbed Diffusion-DPO. We connect this formulation to a multi-step RL approach in the same setting as existing work [6, 11].

We demonstrate the efficacy of Diffusion-DPO by fine-tuning state-of-the-art text-to-image diffusion models, such as Stable Diffusion XL (SDXL)-1.0 [30]. Human evaluators prefer DPO-tuned SDXL images over the SDXL-(base + refinement) model 69% of the time on the PartiPrompts dataset, which represents the state-of-the-art in text-to-image models as measured by human preference. Example generations shown in Fig. 1. Finally, we show that learning from AI feedback (instead of human preferences) using the Diffusion-DPO objective is also effective, a setting where previous works have been unsuccessful [7]. In sum, we introduce a novel paradigm of learning from human preferences for diffusion models and present the resulting state-of-the-art model.

2. Related Work

Aligning Large Language Models LLMs are typically aligned to human preferences using supervised fine-tuning on demonstration data, followed by RLHF. RLHF consists of training a reward function from comparison data on model outputs to represent human preferences and then using reinforcement learning to align the policy model. Prior work [4, 26, 29, 47] has used policy-gradient methods [27, 38] to this end. These methods are successful, but expensive and require extensive hyperparameter tuning [34, 59], and can be prone to reward hacking [10, 12, 41]. Alternative approaches sample base model answers and select based on predicted rewards [3, 5, 14] to use for supervised training [2, 16, 50]. Methods that fine-tune the policy model directly on feedback data [1, 10], or utilize a ranking loss on preference data to directly train the policy model [33, 49, 57, 58] have emerged. The latter set of

methods match RLHF in performance. We build on these fine-tuning methods in this work, specifically, direct preference optimization [33] (DPO). Finally, learning from AI feedback, using pretrained reward models, is promising for efficient scaling of alignment [4, 22].

Aligning Diffusion Models Alignment of diffusion models to human preferences has so far been much less explored than in the case of LLMs. Multiple approaches [30, 36] fine-tune on datasets scored as highly visually appealing by an aesthetics classifier [37], to bias the model to visually appealing generations. Emu [9] finetunes a pretrained model using a small, curated image dataset of high quality photographs with manually written detailed captions to improve visual appeal and text alignment. Other methods [15, 39] recaption existing web-scraped image datasets to improve text fidelity. Caption-aware human preference scoring models are trained on generation preference datasets [21, 52, 55], but the impact of these reward models to the generative space has been limited. DOODL [51] introduces the task of aesthetically improving a single generation iteratively at inference time. DRAFT [7] and AlignProp [31], incorporate a similar approach into training: tuning the generative model to directly increase the reward of generated images. These methods perform well for simple visual appeal criteria, but lack stability and do not work on more nuanced rewards such as text-image alignment from a CLIP model. DPOK and DDPO [6, 11] are RL-based approaches to similarly maximize the scored reward (with distributional constraints) over a relatively limited vocabulary set; the performance of these methods degrades as the number of train/test prompts increases. Diffusion-DPO is unique among alignment approaches in effectively increasing measured human appeal across an open vocabulary (DPOK, DDPO), without increased inference time (DOODL) while maintaining distributional guarantees and improving generic text-image alignment in addition to visual appeal (DRAFT, AlignProp). (see Tab. 1, further discussion in Supp. S1).

Methods	Equal		
	Open Vocab.	Inference Cost	Divergence Control
DPOK[11]	✗	✓	✓
DDPO[6]	✗	✓	✗
DOODL[51]	✓	✗	✗
DRaFT[7], AlignProp[31]	✓	✓	✗
Diffusion-DPO (ours)	✓	✓	✓

Table 1. Method class comparison. Existing methods fail in one or more of: Generalizing to an open vocabulary, maintaining the same inference complexity, avoiding mode collapse/providing distributional guarantees. Diffusion-DPO addresses these issues.

3. Background

3.1. Diffusion Models

Given samples from a data distribution $q(\mathbf{x}_0)$, noise scheduling function α_t and σ_t (as defined in [36]), denoising diffusion models [17, 42, 46] are generative models $p_\theta(\mathbf{x}_0)$ which have a discrete-time reverse process with a Markov structure $p_\theta(\mathbf{x}_{0:T}) = \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ where

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t), \sigma_{t|t-1}^2 \frac{\sigma_{t-1}^2}{\sigma_t^2} \mathbf{I}). \quad (1)$$

Training is performed by minimizing the evidence lower bound (ELBO) associated with this model [20, 45]:

$$L_{\text{DM}} = \mathbb{E}_{\mathbf{x}_0, \epsilon, t, \mathbf{x}_t} [\omega(\lambda_t) \|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|_2^2], \quad (2)$$

with $\epsilon \sim \mathcal{N}(0, \mathbf{I})$, $t \sim \mathcal{U}(0, T)$, $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \alpha_t \mathbf{x}_0, \sigma_t^2 \mathbf{I})$. $\lambda_t = \alpha_t^2 / \sigma_t^2$ is a signal-to-noise ratio [20], $\omega(\lambda_t)$ is a pre-specified weighting function (typically chosen to be constant [17, 44]).

3.2. Direct Preference Optimization

Our approach is an adaption of *Direct Preference Optimization (DPO)* [33], an effective approach for learning from human preference for language models. Abusing notation, we also use \mathbf{x}_0 as random variables for language.

Reward Modeling Estimating human partiality to a generation \mathbf{x}_0 given conditioning \mathbf{c} , is difficult as we do not have access to the latent reward model $r(\mathbf{c}, \mathbf{x}_0)$. In our setting, we assume access only to ranked pairs generated from some conditioning $\mathbf{x}_0^w \succ \mathbf{x}_0^l | \mathbf{c}$, where \mathbf{x}_0^w and \mathbf{x}_0^l denoting the “winning” and “losing” samples. The Bradley-Terry (BT) model stipulates to write human preferences as:

$$p_{\text{BT}}(\mathbf{x}_0^w \succ \mathbf{x}_0^l | \mathbf{c}) = \sigma(r(\mathbf{c}, \mathbf{x}_0^w) - r(\mathbf{c}, \mathbf{x}_0^l)) \quad (3)$$

where σ is the sigmoid function. $r(\mathbf{c}, \mathbf{x}_0)$ can be parameterized by a neural network ϕ and estimated via maximum likelihood training for binary classification:

$$L_{\text{BT}}(\phi) = -\mathbb{E}_{\mathbf{c}, \mathbf{x}_0^w, \mathbf{x}_0^l} [\log \sigma(r_\phi(\mathbf{c}, \mathbf{x}_0^w) - r_\phi(\mathbf{c}, \mathbf{x}_0^l))] \quad (4)$$

where prompt \mathbf{c} and data pairs $\mathbf{x}_0^w, \mathbf{x}_0^l$ are from a static dataset with human-annotated labels.

RLHF RLHF aims to optimize a conditional distribution $p_\theta(\mathbf{x}_0|\mathbf{c})$ (conditioning $\mathbf{c} \sim \mathcal{D}_c$) such that the latent reward model $r(\mathbf{c}, \mathbf{x}_0)$ defined on it is maximized, while regularizing the KL-divergence from a reference distribution p_{ref}

$$\max_{p_\theta} \mathbb{E}_{\mathbf{c} \sim \mathcal{D}_c, \mathbf{x}_0 \sim p_\theta(\mathbf{x}_0|\mathbf{c})} [r(\mathbf{c}, \mathbf{x}_0)] - \beta \mathbb{D}_{\text{KL}} [p_\theta(\mathbf{x}_0|\mathbf{c}) \| p_{\text{ref}}(\mathbf{x}_0|\mathbf{c})] \quad (5)$$

where the hyperparameter β controls regularization.

DPO Objective In Eq. (5) from [33], the unique global optimal solution p_θ^* takes the form:

$$p_\theta^*(\mathbf{x}_0|\mathbf{c}) = p_{\text{ref}}(\mathbf{x}_0|\mathbf{c}) \exp(r(\mathbf{c}, \mathbf{x}_0)/\beta) / Z(\mathbf{c}) \quad (6)$$

where $Z(\mathbf{c}) = \sum_{\mathbf{x}_0} p_{\text{ref}}(\mathbf{x}_0|\mathbf{c}) \exp(r(\mathbf{c}, \mathbf{x}_0)/\beta)$ is the partition function. Hence, the reward function is rewritten as

$$r(\mathbf{c}, \mathbf{x}_0) = \beta \log \frac{p_\theta^*(\mathbf{x}_0|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_0|\mathbf{c})} + \beta \log Z(\mathbf{c}) \quad (7)$$

Using Eq. (4), the reward objective becomes:

$$L_{\text{DPO}}(\theta) = -\mathbb{E}_{\mathbf{c}, \mathbf{x}_0^w, \mathbf{x}_0^l} \left[\log \sigma \left(\beta \log \frac{p_\theta(\mathbf{x}_0^w|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_0^w|\mathbf{c})} - \beta \log \frac{p_\theta(\mathbf{x}_0^l|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_0^l|\mathbf{c})} \right) \right] \quad (8)$$

By this reparameterization, instead of optimizing the reward function r_ϕ and then performing RL, [33] directly optimizes the optimal conditional distribution $p_\theta(\mathbf{x}_0|\mathbf{c})$.

4. DPO for Diffusion Models

In adapting DPO to diffusion models, we consider a setting where we have a fixed dataset $\mathcal{D} = \{(\mathbf{c}, \mathbf{x}_0^w, \mathbf{x}_0^l)\}$ where each example contains a prompt \mathbf{c} and a pairs of images generated from a reference model p_{ref} with human label $\mathbf{x}_0^w \succ \mathbf{x}_0^l$. We aim to learn a new model p_θ which is aligned to the human preferences, with preferred generations to p_{ref} . The primary challenge we face is that the parameterized distribution $p_\theta(\mathbf{x}_0|\mathbf{c})$ is not tractable, as it needs to marginalize out all possible diffusion paths $(\mathbf{x}_1, \dots, \mathbf{x}_T)$ which lead to \mathbf{x}_0 . To overcome this challenge, we utilize the evidence lower bound (ELBO). Here, we introduce latents $\mathbf{x}_{1:T}$ and define $R(\mathbf{c}, \mathbf{x}_{0:T})$ as the reward on the whole chain, such that we can define $r(\mathbf{c}, \mathbf{x}_0)$ as

$$r(\mathbf{c}, \mathbf{x}_0) = \mathbb{E}_{p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{c})} [R(\mathbf{c}, \mathbf{x}_{0:T})]. \quad (9)$$

As for the KL-regularization term in Eq. (5), following prior work [17, 42], we can instead minimize its upper bound joint KL-divergence $\mathbb{D}_{\text{KL}}[p_\theta(\mathbf{x}_{0:T}|\mathbf{c})\|p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c})]$. Plugging this KL-divergence bound and the definition of $r(\mathbf{c}, \mathbf{x}_0)$ (Eq. (9)) back to Eq. (5), we have the objective

$$\max_{p_\theta} \mathbb{E}_{\mathbf{c} \sim \mathcal{D}, \mathbf{x}_{0:T} \sim p_\theta(\mathbf{x}_{0:T}|\mathbf{c})} [r(\mathbf{c}, \mathbf{x}_0)] - \beta \mathbb{D}_{\text{KL}}[p_\theta(\mathbf{x}_{0:T}|\mathbf{c})\|p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c})]. \quad (10)$$

This objective has a parallel formulation as Eq. (5) but defined on path $\mathbf{x}_{0:T}$. It aims to maximize the reward for reverse process $p_\theta(\mathbf{x}_{0:T})$, while matching the distribution of the original reference reverse process. Paralleling Eqs. (6) to (8), this objective can be optimized directly through the

conditional distribution $p_\theta(\mathbf{x}_{0:T})$ via objective:

$$L_{\text{DPO-Diffusion}}(\theta) = -\mathbb{E}_{(\mathbf{x}_0^w, \mathbf{x}_0^l) \sim \mathcal{D}} \log \sigma \left(\beta \mathbb{E}_{\substack{\mathbf{x}_{1:T}^w \sim p_\theta(\mathbf{x}_{1:T}^w|\mathbf{x}_0^w) \\ \mathbf{x}_{1:T}^l \sim p_\theta(\mathbf{x}_{1:T}^l|\mathbf{x}_0^l)}} \left[\log \frac{p_\theta(\mathbf{x}_{0:T}^w)}{p_{\text{ref}}(\mathbf{x}_{0:T}^w)} - \log \frac{p_\theta(\mathbf{x}_{0:T}^l)}{p_{\text{ref}}(\mathbf{x}_{0:T}^l)} \right] \right) \quad (11)$$

We omit \mathbf{c} for compactness (details included in Supp. S2). To optimize Eq. (11), we must sample $\mathbf{x}_{1:T} \sim p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)$. Despite the fact that p_θ contains trainable parameters, this sampling procedure is both (1) *inefficient* as T is usually large ($T = 1000$), and (2) *intractable* since $p_\theta(\mathbf{x}_{1:T})$ represents the reverse process parameterization $p_\theta(\mathbf{x}_{1:T}) = p_\theta(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$. We solve these two issues next.

From Eq. (11), we substitute the reverse decompositions for p_θ and p_{ref} , and utilize Jensen's inequality and the convexity of function $-\log \sigma$ to push the expectation outside. With some simplification, we get the following bound

$$L_{\text{DPO-Diffusion}}(\theta) \leq -\mathbb{E}_{\substack{(\mathbf{x}_0^w, \mathbf{x}_0^l) \sim \mathcal{D}, t \sim \mathcal{U}(0, T), \\ \mathbf{x}_{t-1, t}^w \sim p_\theta(\mathbf{x}_{t-1, t}^w|\mathbf{x}_0^w), \\ \mathbf{x}_{t-1, t}^l \sim p_\theta(\mathbf{x}_{t-1, t}^l|\mathbf{x}_0^l)}} \log \sigma \left(\beta T \log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \beta T \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right) \quad (12)$$

Efficient training via gradient descent is now possible. However, sampling from reverse joint $p_\theta(\mathbf{x}_{t-1}, \mathbf{x}_t|\mathbf{x}_0, \mathbf{c})$ is still intractable and r of Eq. (9) has an expectation over $p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)$. So we approximate the reverse process $p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)$ with the forward $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$ (an alternative scheme in Supp. S2). With some algebra, this yields:

$$L(\theta) = -\mathbb{E}_{(\mathbf{x}_0^w, \mathbf{x}_0^l) \sim \mathcal{D}, t \sim \mathcal{U}(0, T), \mathbf{x}_t^w \sim q(\mathbf{x}_t^w|\mathbf{x}_0^w), \mathbf{x}_t^l \sim q(\mathbf{x}_t^l|\mathbf{x}_0^l)} \log \sigma \left(-\beta T \left(\mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1}^w|\mathbf{x}_0^w)\|p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)) - \mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1}^w|\mathbf{x}_0^w)\|p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)) - \mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1}^l|\mathbf{x}_0^l)\|p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)) + \mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1}^l|\mathbf{x}_0^l)\|p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)) \right) \right) \quad (13)$$

Using Eq. (1) and algebra, the above loss simplifies to:

$$L(\theta) = -\mathbb{E}_{(\mathbf{x}_0^w, \mathbf{x}_0^l) \sim \mathcal{D}, t \sim \mathcal{U}(0, T), \mathbf{x}_t^w \sim q(\mathbf{x}_t^w|\mathbf{x}_0^w), \mathbf{x}_t^l \sim q(\mathbf{x}_t^l|\mathbf{x}_0^l)} \log \sigma \left(-\beta T \omega(\lambda_t) \left(\|\boldsymbol{\epsilon}^w - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t^w, t)\|_2^2 - \|\boldsymbol{\epsilon}^w - \boldsymbol{\epsilon}_{\text{ref}}(\mathbf{x}_t^w, t)\|_2^2 - (\|\boldsymbol{\epsilon}^l - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t^l, t)\|_2^2 - \|\boldsymbol{\epsilon}^l - \boldsymbol{\epsilon}_{\text{ref}}(\mathbf{x}_t^l, t)\|_2^2) \right) \right) \quad (14)$$

where $\mathbf{x}_t^* = \alpha_t \mathbf{x}_0^* + \sigma_t \boldsymbol{\epsilon}^*$, $\boldsymbol{\epsilon}^* \sim \mathcal{N}(0, I)$ is a draw from $q(\mathbf{x}_t^*|\mathbf{x}_0^*)$ (Eq. (2)). $\lambda_t = \alpha_t^2/\sigma_t^2$ is the signal-to-noise ratio,

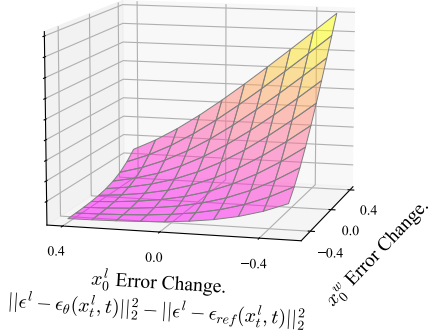


Figure 2. Loss surface visualization. Loss can be decreased by improving at denoising x_0^w and worsening for x_0^l . A larger β increases surface curvature.

$\omega(\lambda_t)$ a weighting function (constant in practice [17, 20]). We factor the constant T into β . This loss encourages ϵ_θ to improve more at denoising x_t^w than x_t^l , visualization in Fig. 2. We also derive Eq. (14) as a multi-step RL approach in the same setting as DDPO and DPOK [6, 11] (Supp. S3) but as an off-policy algorithm, which justifies our sampling choice in Eq. 13. A noisy preference model perspective yields the same objective (Supp. S4).

5. Experiments

5.1. Setting

Models and Dataset: We demonstrate the efficacy of Diffusion-DPO across a range of experiments. We use the objective from Eq. (14) to fine-tune Stable Diffusion 1.5 (SD1.5) [36] and the state-of-the-art open-source model Stable Diffusion XL-1.0 (SDXL) [30] base model. We train on the **Pick-a-Pic** [21] dataset, which consists of pairwise preferences for images generated by SDXL-beta and Dreamlike, a fine-tuned version of SD1.5. The prompts and preferences were collected from users of the Pick-a-Pic web application (see [21] for details). We use the larger Pick-a-Pic v2 dataset. After excluding the $\sim 12\%$ of pairs with ties, we end up with 851,293 pairs, with 58,960 unique prompts.

Hyperparameters We use AdamW [24] for SD1.5 experiments, and Adafactor [40] for SDXL to save memory. An effective batch size of 2048 (pairs) is used; training on 16 NVIDIA A100 GPUs with a local batch size of 1 pair and gradient accumulation of 128 steps. We train at fixed square resolutions. A learning rate of $\frac{2000}{\beta} 2.048 \cdot 10^{-8}$ is used with 25% linear warmup. The inverse scaling is motivated by the norm of the DPO objective gradient being proportional to β (the divergence penalty parameter) [33]. For both SD1.5 and SDXL, we find $\beta \in [2000, 5000]$ to offer good performance (Supp. S5). We present main SD1.5 results with $\beta = 2000$ and SDXL results with $\beta = 5000$.

Evaluation We automatically validate checkpoints with the 500 unique prompts of the Pick-a-Pic validation set: measuring median PickScore reward of generated images. PickScore [21] is a caption-aware scoring model trained on Pick-a-Pic (v1) to estimate human-perceived image quality. For final testing, we generate images using the baseline and Diffusion-DPO-tuned models conditioned on captions from the Partiprompt [56] and HPSv2 [52] benchmarks (1632 and 3200 captions respectively). While DDPO [6] is a related method, we did not observe stable improvement when training from public implementations on Pick-a-Pic. We employ labelers on Amazon Mechanical Turk to compare generations under three different criteria: Q1 General Preference (*Which image do you prefer given the prompt?*), Q2 Visual Appeal (prompt not considered) (*Which image is more visually appealing?*) Q3 Prompt Alignment (*Which image better fits the text description?*). Five responses are collected for each comparison with majority vote (3+) being considered the collective decision.

5.2. Primary Results: Aligning Diffusion Models

First, we show that the outputs of the Diffusion-DPO-finetuned SDXL model are significantly preferred over the baseline SDXL-base model. In the Partiprompt evaluation (Fig. 3-top left), DPO-SDXL is preferred 70.0% of the time for General Preference (Q1), and obtains a similar win-rate in assessments of both Visual Appeal (Q2) and Prompt Alignment (Q3). Evaluation on the HPS benchmark (Fig. 3-top right) shows a similar trend, with a General Preference win rate of 64.7%. We also score the DPO-SDXL HPSv2 generations with the HPSv2 reward model, achieving an average reward of 28.16, topping the leaderboard [53].

We display qualitative comparisons to SDXL-base in Fig. 3 (bottom). Diffusion-DPO produces more appealing imagery, with vivid arrays of colors, dramatic lighting, good composition, and realistic people/animal anatomy. While all SDXL images satisfy the prompting criteria to some degree, the DPO generations appear superior, as confirmed by the crowdsourced study. We do note that preferences are not universal, and while the most common shared preference is towards energetic and dramatic imagery, others may prefer quieter/subtler scenes. The area of personal or group preference tuning is an exciting area of future work.

After this parameter-equal comparison with SDXL-base, we compare SDXL-DPO to the complete SDXL pipeline, consisting of the base model and the refinement model (Fig. 4). The refinement model is an image-to-image diffusion model that improves visual quality of generations, and is especially effective on detailed backgrounds and faces. In our experiments with PartiPrompts and HPSv2, SDXL-DPO (3.5B parameters, SDXL-base architecture only), handily beats the complete SDXL model (6.6B parameters). In the General Preference question, it has a benchmark win

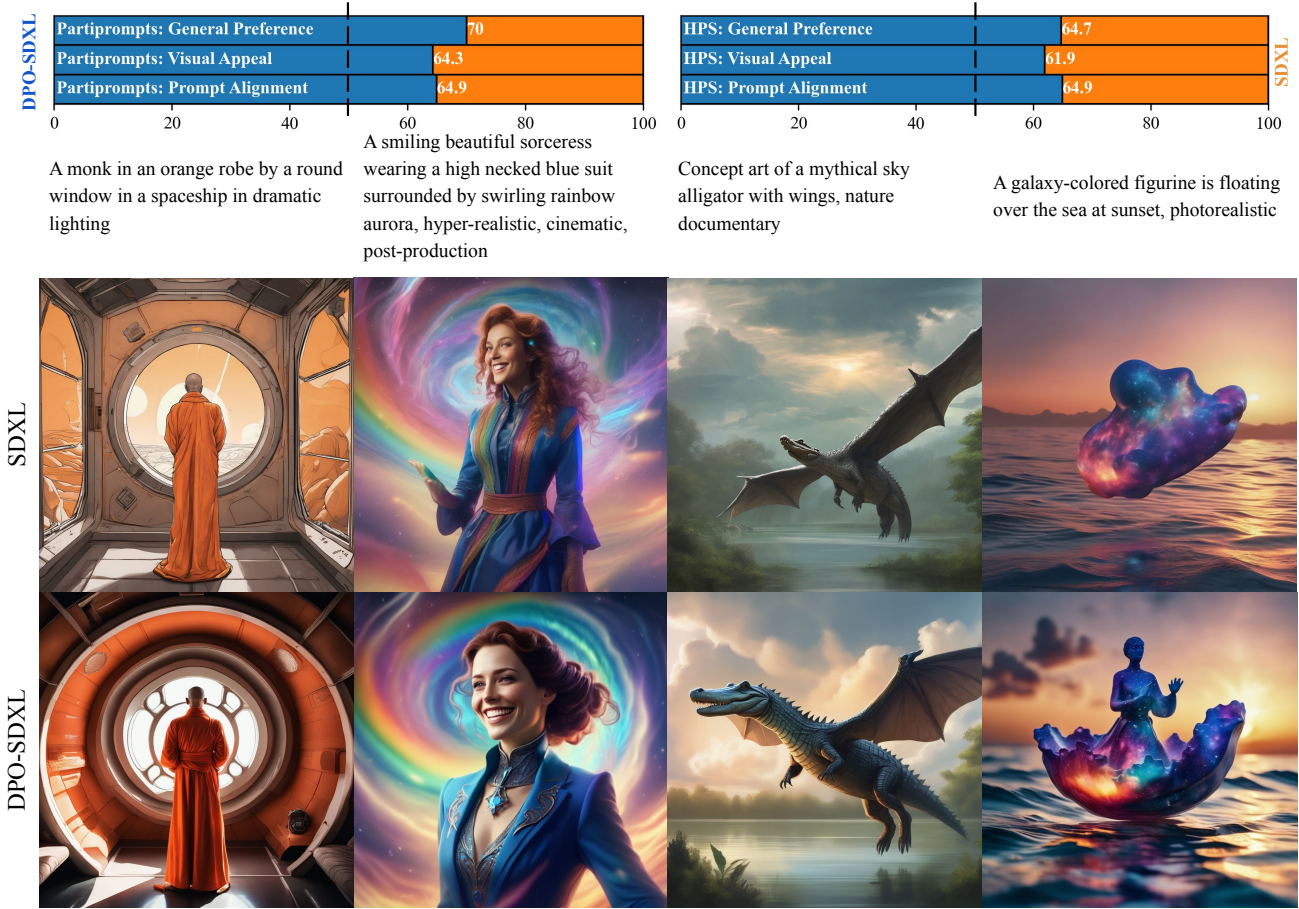


Figure 3. (Top) **DPO-SDXL** significantly outperforms **SDXL** in human evaluation. (L) PartiPrompts and (R) HPSv2 benchmark results across three evaluation questions, majority vote of 5 labelers. (Bottom) Qualitative comparisons between **SDXL** and **DPO-SDXL**. **DPO-SDXL** demonstrates superior prompt following and realism. **DPO-SDXL** outputs are better aligned with human aesthetic preferences, favoring high contrast, vivid colors, fine detail, and focused composition. They also capture fine-grained textual details more faithfully.

rate of 69% and 64% respectively, comparable to its win rate over **SDXL**-base alone. This is explained by the ability of the **DPO**-tuned model (Fig. 4, bottom) to generate fine-grained details and its strong performance across different image categories. While the refinement model is especially good at improving the generation of human details, the win rate of **Diffusion-DPO** on the *People* category in **PartiPrompt** dataset over the base + refiner model is still an impressive 67.2% (compared to 73.4% over the base).

5.3. Image-to-Image Editing

Image-to-image translation performance also improves after **Diffusion-DPO** tuning. We test **DPO-SDXL** on **TEd-Bench** [19], a text-based image-editing benchmark of 100 real image-text pairs, using **SDEdit** [25] with noise strength 0.6. Labelers are shown the original image and **SDXL**/**DPO-SDXL** edits and asked “Which edit do you prefer given the text?” **DPO-SDXL** is preferred 65% of the

time, **SDXL** 24%, with 11% draws. We show qualitative **SDEdit** results on color layouts (strength 0.98) in Fig. 5.

5.4. Learning from AI Feedback

In LLMs, learning from AI feedback has emerged as a strong alternative to learning from human preferences [22]. **Diffusion-DPO** can also admit learning from AI feedback by directly ranking generated pairs into (y_w, y_l) using a pretrained scoring network. We use **HPSv2** [52] for an alternate prompt-aware human preference estimate, **CLIP** (OpenCLIP ViT-H/14) [18, 32] for text-image alignment, **Aesthetic Predictor** [37] for non-text-based visual appeal, and **PickScore**. We run all experiments on **SD 1.5** with $\beta = 5000$ for 1000 steps. Training on the prompt-aware **PickScore** and **HPS** preference estimates increase the win rate for both raw visual appeal and prompt alignment (Fig. 6). We note that **PickScore** feedback is interpretable as pseudo-labeling the **Pick-a-Pic** dataset—a form of data

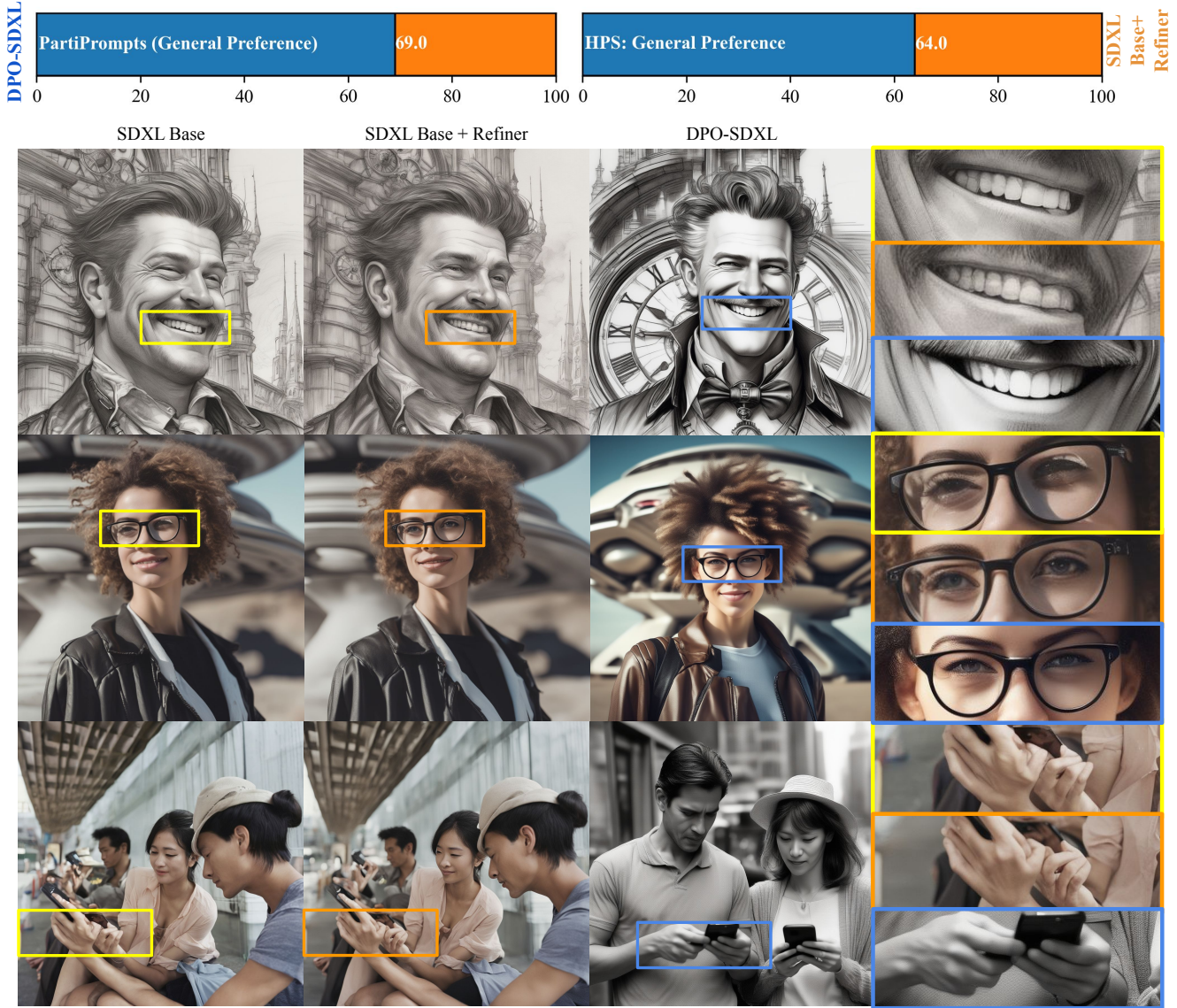


Figure 4. DPO-SDXL (base only) significantly outperforms the much larger SDXL-(base+refinement) model pipeline in human evaluations on the PartiPrompts and HPS datasets. While the SDXL refinement model is used to touch up details from the output of SDXL-base, the ability to generate high quality details has been naturally distilled into DPO-SDXL by human preference. Among other advantages, DPO-SDXL shows superior generation of anatomical features such as teeth, hands, and eyes. Prompts: *close up headshot, steampunk middle-aged man, slick hair big grin in front of gigantic clocktower, pencil sketch / close up headshot, futuristic young woman with glasses, wild hair sly smile in front of gigantic UFO, dslr, sharp focus, dynamic composition / A man and woman using their cellphones, photograph*

cleaning [54, 60]. Training for Aesthetics and CLIP improves those capabilities more specifically, in the case of Aesthetics at the expense of CLIP. The ability to train for text-image alignment via CLIP is a noted improvement over prior work [7]. Moreover, training SD1.5 on the pseudo-labeled PickScore dataset ($\beta = 5000$, 2000 steps) outperforms training on the raw labels. On the General Preference Partiprompt question, the win-rate of DPO increases from 59.8% to 63.3%, indicating that learning from AI feedback

can be a promising direction for diffusion model alignment.

5.5. Analysis

Implicit Reward Model As a consequence of the theoretical framework, our DPO scheme implicitly learns a reward model and can estimate the differences in rewards between two images by taking an expectation over the inner term of Eq. (14) (details in Supp. S4.1). We estimate over 10 random $t \sim \mathcal{U}\{0, 1\}$. Our learned models (DPO-SD1.5

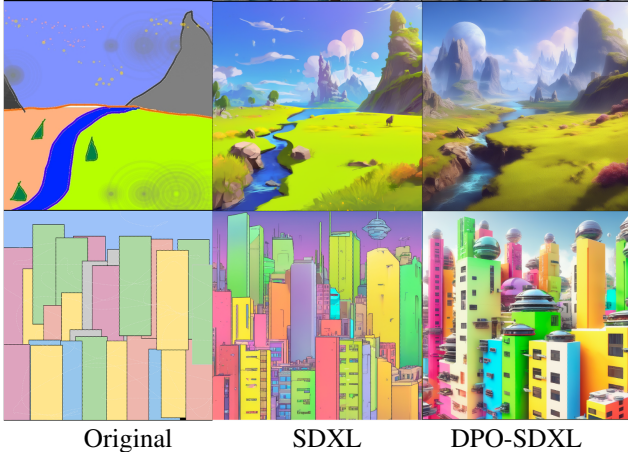


Figure 5. Diffusion-DPO generates more visually appealing images in the downstream image-to-image translation task. Comparisons of using SDEdit [25] from color layouts. Prompts are "A fantasy landscape, trending on artstation" (top), "High-resolution rendering of a crowded colorful sci-fi city" (bottom).

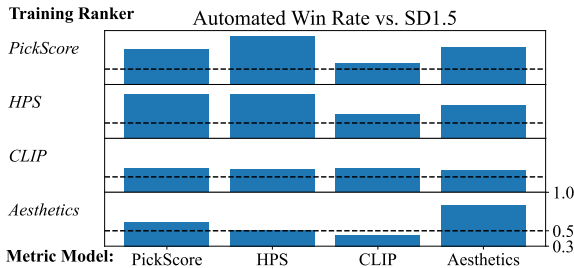


Figure 6. Automated head-to-head win rates under reward models (x labels, columns) for SD1.5 DPO-tuned on the "preferences" of varied scoring networks (y labels, rows). Example: Tuning on *Aesthetics* preferences (bottom row) achieves high *Aesthetics* scores but has lower text-image alignment as measured by CLIP.

Model	PS	HPS	CLIP	Aes.	DPO-SD1.5	DPO-SDXL
Acc.	64.2	59.3	57.1	51.4	60.8	72.0

Table 2. Preference accuracy on the Pick-a-Pic (v2) validation set. The v1-trained PickScore has seen the evaluated data.

and DPO-SDXL) perform well at binary preference classification (Tab. 2), with DPO-SDXL exceeding all existing recognition models on this split. These results shows that the implicit reward parameterization in the Diffusion-DPO objective has comprable expressivity and generalization as the classical reward modelling objective/architecture.

Training Data Quality Fig. 7 shows that despite SDXL being superior to the training data (including the y_w), as measured by Pickscore, DPO training improves its perfor-

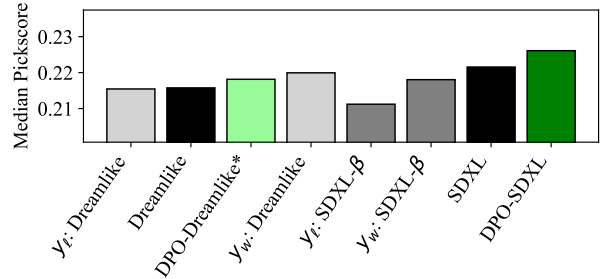


Figure 7. Diffusion-DPO improves on the baseline Dreamlike and SDXL models, when finetuned on both in-distribution data (in case of Dreamlike) and out-of-distribution data (in case of SDXL). y_t and y_w denote the Pickscore of winning and losing samples.

mance substantially. In this experiment, we confirm that Diffusion-DPO can improve on in-distribution preferences as well, by training ($\beta = 5k$, 2000 steps) the Dreamlike model on a subset of the Pick-a-Pic dataset generated by the Dreamlike model alone. This subset represents 15% of the original dataset. Dreamlike-DPO improves on the baseline model, though the performance improvement is limited, perhaps because of the small size of the dataset.

Supervised Fine-tuning (SFT) is beneficial in the LLM setting as initial pretraining prior to preference training. To evaluate SFT in our setting, we fine-tune models on the preferred (x, y_w) pairs of the Pick-a-Pic dataset. We train for the same length schedule as DPO using a learning rate of $1e - 9$ and observe convergence. While SFT improves vanilla SD1.5 (55.5% win rate over base model), any amount of SFT deteriorates the performance of SDXL, even at lower learning rates. This contrast is attributable to the much higher quality of Pick-a-Pic generations vs. SD1.5, as they are obtained from SDXL-beta and Dreamlike. In contrast, the SDXL-1.0 base model is superior to the Pick-a-Pic dataset models. See Supp. S6 for further discussion.

6. Conclusion

In this work, we introduce Diffusion-DPO: a method that enables diffusion models to directly learn from human feedback in an open-vocabulary setting for the first time. We fine-tune SDXL-1.0 using the Diffusion-DPO objective and the Pick-a-Pic (v2) dataset to create a new state-of-the-art for open-source text-to-image generation models as measured by generic preference, visual appeal, and prompt alignment. We additionally demonstrate that DPO-SDXL outperforms even the SDXL base plus refinement model pipeline, despite only employing 53% of the total model parameters. Dataset cleaning/scaling is a promising future direction as we observe preliminary data cleaning improving performance (Sec. 5.4). While DPO-Diffusion is an of-

fine algorithm, we anticipate online learning methods to be another driver of future performance. There are also exciting application variants such as tuning to the preferences of individuals or small groups.

Ethics The performance of Diffusion-DPO is impressive, but any effort in text-to-image generation presents ethical risks, particularly when data are web-collected. Generations of harmful, hateful, fake or sexually explicit content are known risk vectors. Beyond that, this approach is increasingly subject to the biases of the participating labelers (in addition to the biases present in the pretrained model); Diffusion-DPO can learn and propagate these preferences. As a result, a diverse and representative set of labelers is essential – whose preferences in turn become encoded in the dataset. Furthermore, a portion of user-generated Pick-a-Pic prompts are overtly sexual, and even innocuous prompts may deliver images that skew more suggestively (particularly for prompts that hyper-sexualize women). Finally, as with all text-to-image models, the image produced will not always match the prompt. Hearteningly though, some of these scenarios can be addressed at a dataset level, and data filtering is also possible. Regardless, we will not open source nor otherwise make available our model until we add additional safety filtering to ensure that toxic content is remediated.

References

- [1] Model index for researchers, 2023. [2](#)
- [2] Thomas Anthony, Zheng Tian, and David Barber. Thinking fast and slow with deep learning and tree search. *Neural Information Processing Systems*, 2017. [2](#)
- [3] Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, and Jared Kaplan. A general language assistant as a laboratory for alignment, 2021. [2](#)
- [4] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai: Harmlessness from ai feedback, 2022. [2](#), [3](#)
- [5] Michiel A. Bakker, Martin J. Chadwick, Hannah R. Sheahan, Michael Henry Tessler, Lucy Campbell-Gillingham, Jan Balaguer, Nat McAleese, Amelia Glaese, John Aslanides, Matthew M. Botvinick, and Christopher Summerfield. Fine-tuning language models to find agreement among humans with diverse preferences. *Neural Information processing systems*, 2022. [2](#)
- [6] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023. [1](#), [2](#), [3](#), [5](#), [4](#)
- [7] Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models on differentiable rewards, 2023. [1](#), [2](#), [3](#), [7](#)
- [8] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. Vqgan-clip: Open domain image generation and editing with natural language guidance, 2022. [1](#)
- [9] Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiang Wang, Rui Wang, Peizhao Zhang, Simon Vandenhende, Xiaofang Wang, Abhimanyu Dubey, et al. Emu: Enhancing image generation models using photogenic needles in a haystack. *arXiv preprint arXiv:2309.15807*, 2023. [1](#), [3](#), [2](#)
- [10] Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback, 2023. [2](#)
- [11] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpoc: Reinforcement learning for fine-tuning text-to-image diffusion models. *arXiv preprint arXiv:2305.16381*, 2023. [1](#), [2](#), [3](#), [5](#), [4](#)
- [12] Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization, 2022. [2](#)
- [13] Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, Matthieu Geist, and Stefano Ermon. Iq-learn: Inverse soft-q learning for imitation. *Neural Information Processing Systems*, 2021. [5](#)
- [14] Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. Improving alignment of dialogue agents via targeted human judgements, 2022. [2](#)
- [15] Gabriel Goh, James Betker, Li Jing, Aditya Ramesh, Tim Brooks, Jianfeng Wang, Lindsey Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Prafulla Dhariwal, Casey Chu, Joy Jiao, Jong Wook Kim, Alex Nichol, Yang Song, Lijuan Wang, and Tao Xu. Improving image generation with better captions. 2023. [3](#), [10](#), [12](#)

- [16] Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, Wolfgang Macherey, Arnaud Doucet, Orhan Firat, and Nando de Freitas. Reinforced self-training (rest) for language modeling, 2023. **2**
- [17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. pages 6840–6851, 2020. **3, 4, 5**
- [18] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hananeh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. If you use this software, please cite it as below. **6**
- [19] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6007–6017, 2023. **6**
- [20] Diederik P. Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. 2021. **3, 5, 4**
- [21] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *arXiv preprint arXiv:2305.01569*, 2023. **3, 5, 8**
- [22] Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. Rlaif: Scaling reinforcement learning from human feedback with ai feedback, 2023. **3, 6**
- [23] Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review, 2018. **4**
- [24] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. **5**
- [25] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. **6, 8**
- [26] Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, and Nat McAleese. Teaching language models to support answers with verified quotes, 2022. **2**
- [27] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning, 2016. **2**
- [28] OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023. **1**
- [29] 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, 2022. **2**
- [30] 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. **2, 3, 5, 1**
- [31] Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023. **1, 3**
- [32] 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*, pages 8748–8763. PMLR, 2021. **6, 1**
- [33] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2023. **1, 2, 3, 4, 5, 8**
- [34] Rajkumar Ramamurthy, Prithviraj Ammanabrolu, Kianté Brantley, Jack Hessel, Rafet Sifa, Christian Bauckhage, Hananeh Hajishirzi, and Yejin Choi. Is reinforcement learning (not) for natural language processing: Benchmarks, baselines, and building blocks for natural language policy optimization, 2022. **2**
- [35] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Stable diffusion 2. <https://huggingface.co/stabilityai/stable-diffusion-2>. Accessed: 2023 - 11 - 16. **1**
- [36] 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 (CVPR)*, pages 10684–10695, 2022. **1, 3, 5**
- [37] Christoph Schuhmann. Laion-aesthetics. <https://laion.ai/blog/laion-aesthetics/>, 2022. Accessed: 2023 - 11 - 10. **3, 6, 1**
- [38] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. **2**
- [39] Eyal Segalis, Dani Valevski, Danny Lumen, Yossi Matias, and Yaniv Leviathan. A picture is worth a thousand words: Principled recaptioning improves image generation, 2023. **3**
- [40] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR, 2018. **5**
- [41] Joar Skalse, Nikolaus H. R. Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and characterizing reward hacking. *Neural Information Processing Systems*, 2022. **2**
- [42] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265, 2015. **3, 4**
- [43] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. **6**
- [44] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019. **3**

- [45] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. In *Neural Information Processing Systems*, 2021. 3
- [46] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021. 3
- [47] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. *Neural Information Processing Systems*, 18, 2020. 2
- [48] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. 1
- [49] Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of Lm alignment, 2023. 2
- [50] Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process- and outcome-based feedback, 2022. 2
- [51] Bram Wallace, Akash Gokul, Stefano Ermon, and Nikhil Naik. End-to-end diffusion latent optimization improves classifier guidance. *arXiv preprint arXiv:2303.13703*, 2023. 3, 1
- [52] 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. *arXiv preprint arXiv:2306.09341*, 2023. 3, 5, 6
- [53] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Hpsv2 github. <https://github.com/tgxs002/HPSv2/tree/master>, 2023. Accessed: 2023 - 11 - 15. 5
- [54] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10687–10698, 2020. 7
- [55] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imageward: Learning and evaluating human preferences for text-to-image generation. *arXiv preprint arXiv:2304.05977*, 2023. 3
- [56] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation, 2022. 5
- [57] Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf: Rank responses to align language models with human feedback without tears. *Neural Information Processing Systems*, 2023. 2
- [58] Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J. Liu. Slic-hf: Sequence likelihood calibration with human feedback, 2023. 2
- [59] Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Qin Liu, Yuhao Zhou, Limao Xiong, Lu Chen, Zhiheng Xi, Nuo Xu, Wenbin Lai, Minghao Zhu, Cheng Chang, Zhangyue Yin, Rongxiang Weng, Wensen Cheng, Haoran Huang, Tianxiang Sun, Hang Yan, Tao Gui, Qi Zhang, Xipeng Qiu, and Xuanjing Huang. Secrets of rlhf in large language models part i: Ppo, 2023. 2
- [60] Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin Dogus Cubuk, and Quoc Le. Rethinking pre-training and self-training. *Advances in neural information processing systems*, 33:3833–3845, 2020. 7

Diffusion Model Alignment Using Direct Preference Optimization

Supplementary Material

S1. Comparisons to existing work

RL-Based Methods such as [6, 11] have shown effectiveness in operating on a limited set of prompts (< 10 and < 1000 respectively) but do not generalize as well to the open-vocabulary setting as shown in [7, 31]. We found this in our experiments as well, where training using the DDPO scheme did not improve PickScore over the baseline model over a sweep of hyperparameters.

While DDPO [6] is an RL-based method as is DPOK [11], their target objective and distributional guarantees are different. Specifically, DDPO purely aims to optimize the reward function without any KL-regularization

$$\mathbb{E}_{\mathbf{c} \sim p(\mathbf{c}), \mathbf{x}_0 \sim p_\theta(\mathbf{x}_0 | \mathbf{c})} r(\mathbf{x}_0, \mathbf{c}) \quad (15)$$

while DPOK adds in a term governing KL-regularization between the learned distribution and a reference distribution as in our setting. This means that DDPO is optimizing the same objective as DRaFT and AlignProp ([7, 31]) but via RL instead of gradient descent through the diffusion chain. DDPO uses early stopping in lieu of distributional control.

Additionally, through the score function policy gradient estimator employed by DDPO it is observable why the method struggles with open vocabulary generation. The gradient estimation used is

$$\nabla_\theta \mathcal{J}_{\text{DDRL}} = \mathbb{E} \sum_{t=0}^T \frac{p_\theta(x_{t-1} | x_t, \mathbf{c})}{p_{\theta_{\text{old}}}(x_{t-1} | x_t, \mathbf{c})} \nabla_\theta \log p_\theta(x_{t-1} | x_t, \mathbf{c}) r(x_0, \mathbf{c}) \quad (16)$$

Here the trajectories $\{\mathbf{x}_T, \mathbf{x}_{T-1}, \dots, \mathbf{x}_0\}$ are generated by the *original* model $p_{\theta_{\text{old}}}$. In this formulation, the term $\frac{p_\theta(x_{t-1} | x_t, \mathbf{c})}{p_{\theta_{\text{old}}}(x_{t-1} | x_t, \mathbf{c})}$ simply is an importance weighter which scales gradient contributions based on the relevance of the sample (as determined by how aligned the learned and reference model predictions are). Since the trajectories are generated by the “old” (reference) model, $r(x_0, \mathbf{c})$ is only a weighting in the latter term $\nabla_\theta \log p_\theta(x_{t-1} | x_t, \mathbf{c}) r(x_0, \mathbf{c})$. The gradient encourages higher likelihoods for generations of high reward, but makes no distinction about the diversity of those generations. High-reward prompts can dominate the gradient trajectory, while generations considered lower-reward are ignored or discouraged. This stands in contrast to the DPO framework where the likelihood of a generation is contrasted against another with the same conditioning. This normalization across conditioning prevents sets of \mathbf{c} being considered unimportant/undesirable and not being optimized for. In Diffusion-DPO, conditionings with all types of reward magnitudes are weighted equally towards the \mathbf{x}_0^w and away from the \mathbf{x}_0^l .

Inference Time-Optimization namely DOODL [51], does not learn any new model parameters, instead optimizing diffusion latents to improve some criterion on the generated image similar to CLIP+VQGAN[8]. This runtime compute increases inference cost by more than an order of magnitude.

Reward Maximization Training such as [7, 31] amortize the cost of DOODL from runtime to training. They train by generating images from text prompts, computing a reward loss on the images, and backpropagating gradients through the generative process to improve the loss. While effective in the open-vocabulary setting (also training on Pick-a-Pic prompts), these methods provide no distributional guarantees (unlike the control via β in Diffusion-DPO) and suffer from mode collapse with over-training. These methods do not generalize to all reward functions, with [7] noting the inability of DRaFT to improve image-text alignment using CLIP[32] as a reward function. In contrast, Diffusion-DPO can improve image-text alignment using CLIP preference, as shown in Sec. 5.4. Furthermore, only differentiable rewards can be optimized towards in the reward maximization setting. This necessitates not only data collection but also reward model training.

Dataset Curation As discussed, models such as StableDiffusion variants [30, 36] train on laion-aesthetics [37] to bias the model towards more visually appealing outputs. Concurrent work Emu [9] takes this approach to an extreme. Instead of training on any images from a web-scale dataset which pass a certain model score threshold, they employ a multi-stage

pipeline where such filtering is only the first stage. Subsequently, crowd workers filter the subset down using human judgement and at the final stage expert in photography are employed to create the dataset. While effective, this process has several drawbacks compared to Diffusion-DPO. First, necessitating training on existing data can be a bottleneck, both in terms of scale and potential applications. While [9] reports lesser text faithfulness improvements as well, these are likely due to the hand-written captions, a much more costly data collection stage than preferences. The Emu pipeline is not generalizable to different types of feedback as DPO is (e.g. outside of recaptioning it is non-obvious how such an approach can improve text-image alignment).

S2. Details of the Primary Derivation

Starting from Eq. (5), we have

$$\begin{aligned}
& \min_{p_\theta} -\mathbb{E}_{p_\theta(\mathbf{x}_0|\mathbf{c})} r(\mathbf{c}, \mathbf{x}_0)/\beta + \mathbb{D}_{\text{KL}}(p_\theta(\mathbf{x}_0|\mathbf{c})||p_{\text{ref}}(\mathbf{x}_0|\mathbf{c})) \\
& \leq \min_{p_\theta} -\mathbb{E}_{p_\theta(\mathbf{x}_0|\mathbf{c})} r(\mathbf{c}, \mathbf{x}_0)/\beta + \mathbb{D}_{\text{KL}}(p_\theta(\mathbf{x}_{0:T}|\mathbf{c})||p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c})) \\
& = \min_{p_\theta} -\mathbb{E}_{p_\theta(\mathbf{x}_{0:T}|\mathbf{c})} R(\mathbf{c}, \mathbf{x}_{0:T})/\beta + \mathbb{D}_{\text{KL}}(p_\theta(\mathbf{x}_{0:T}|\mathbf{c})||p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c})) \\
& = \min_{p_\theta} \mathbb{E}_{p_\theta(\mathbf{x}_{0:T}|\mathbf{c})} \left(\log \frac{p_\theta(\mathbf{x}_{0:T}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c}) \exp(R(\mathbf{c}, \mathbf{x}_{0:T})/\beta)/Z(\mathbf{c})} - \log Z(\mathbf{c}) \right) \\
& = \min_{p_\theta} \mathbb{D}_{\text{KL}}(p_\theta(\mathbf{x}_{0:T}|\mathbf{c})||p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c}) \exp(R(\mathbf{c}, \mathbf{x}_{0:T})/\beta)/Z(\mathbf{c})).
\end{aligned} \tag{17}$$

where $Z(\mathbf{c}) = \sum_{\mathbf{x}} p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c}) \exp(r(\mathbf{c}, \mathbf{x}_0)/\beta)$ is the partition function. The optimal $p_\theta^*(\mathbf{x}_{0:T}|\mathbf{c})$ of Equation (17) has a unique closed-form solution:

$$p_\theta^*(\mathbf{x}_{0:T}|\mathbf{c}) = p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c}) \exp(R(\mathbf{c}, \mathbf{x}_{0:T})/\beta)/Z(\mathbf{c}),$$

Therefore, we have the reparameterization of reward function

$$R(\mathbf{c}, \mathbf{x}_{0:T}) = \beta \log \frac{p_\theta^*(\mathbf{x}_{0:T}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c})} + \beta \log Z(\mathbf{c}).$$

Plug this into the definition of r , hence we have

$$r(\mathbf{c}, \mathbf{x}_0) = \beta \mathbb{E}_{p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{c})} \left[\log \frac{p_\theta^*(\mathbf{x}_{0:T}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_{0:T}|\mathbf{c})} \right] + \beta \log Z(\mathbf{c}).$$

Substituting this reward reparameterization into maximum likelihood objective of the Bradley-Terry model as Eq. (4), the partition function cancels for image pairs, and we get a maximum likelihood objective defined on diffusion models, its per-example formula is:

$$L_{\text{DPO-Diffusion}}(\theta) = -\log \sigma \left(\beta \mathbb{E}_{\mathbf{x}_{1:T}^w \sim p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0^w), \mathbf{x}_{1:T}^l \sim p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0^l)} \left[\log \frac{p_\theta(\mathbf{x}_{0:T}^w)}{p_{\text{ref}}(\mathbf{x}_{0:T}^w)} - \log \frac{p_\theta(\mathbf{x}_{0:T}^l)}{p_{\text{ref}}(\mathbf{x}_{0:T}^l)} \right] \right)$$

where $\mathbf{x}_0^w, \mathbf{x}_0^l$ are from static dataset, we drop \mathbf{c} for simplicity.

An approximation for reverse process Since sampling from $p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)$ is intractable, we utilize $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$ for approximation.

$$\begin{aligned}
L_1(\theta) &= -\log \sigma \left(\beta \mathbb{E}_{\mathbf{x}_{1:T}^w \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0^w), \mathbf{x}_{1:T}^l \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0^l)} \left[\log \frac{p_\theta(\mathbf{x}_{0:T}^w)}{p_{\text{ref}}(\mathbf{x}_{0:T}^w)} - \log \frac{p_\theta(\mathbf{x}_{0:T}^l)}{p_{\text{ref}}(\mathbf{x}_{0:T}^l)} \right] \right) \\
&= -\log \sigma \left(\beta \mathbb{E}_{\mathbf{x}_{1:T}^w \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0^w), \mathbf{x}_{1:T}^l \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0^l)} \left[\sum_{t=1}^T \log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right) \\
&= -\log \sigma \left(\beta \mathbb{E}_{\mathbf{x}_{1:T}^w \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0^w), \mathbf{x}_{1:T}^l \sim q(\mathbf{x}_{1:T}|\mathbf{x}_0^l)} T \mathbb{E}_t \left[\log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right) \\
&= -\log \sigma \left(\beta T \mathbb{E}_t \mathbb{E}_{\mathbf{x}_{t-1,t}^w \sim q(\mathbf{x}_{t-1,t}|\mathbf{x}_0^w), \mathbf{x}_{t-1,t}^l \sim q(\mathbf{x}_{t-1,t}|\mathbf{x}_0^l)} \left[\log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right) \\
&= -\log \sigma \left(\beta T \mathbb{E}_t \mathbb{E}_{\mathbf{x}_t^w \sim q(\mathbf{x}_t|\mathbf{x}_0^w), \mathbf{x}_t^l \sim q(\mathbf{x}_t|\mathbf{x}_0^l)} \right. \\
&\quad \left. \mathbb{E}_{\mathbf{x}_{t-1}^w \sim q(\mathbf{x}_{t-1}|\mathbf{x}_0^w), \mathbf{x}_{t-1}^l \sim q(\mathbf{x}_{t-1}|\mathbf{x}_0^l)} \left[\log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right)
\end{aligned} \tag{18}$$

By Jensen's inequality, we have

$$\begin{aligned}
L_1(\theta) &\leq -\mathbb{E}_{t, \mathbf{x}_t^w \sim q(\mathbf{x}_t|\mathbf{x}_0^w), \mathbf{x}_t^l \sim q(\mathbf{x}_t|\mathbf{x}_0^l)} \log \sigma \left(\right. \\
&\quad \left. \beta T \mathbb{E}_{\mathbf{x}_{t-1}^w \sim q(\mathbf{x}_{t-1}|\mathbf{x}_0^w), \mathbf{x}_{t-1}^l \sim q(\mathbf{x}_{t-1}|\mathbf{x}_0^l)} \left[\log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right) \\
&= -\mathbb{E}_{t, \mathbf{x}_t^w \sim q(\mathbf{x}_t|\mathbf{x}_0^w), \mathbf{x}_t^l \sim q(\mathbf{x}_t|\mathbf{x}_0^l)} \log \sigma \left(-\beta T (\mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1}^w|\mathbf{x}_0^w, t) \| p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)) - \mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1}^w|\mathbf{x}_0^w, t) \| p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w))) \right. \\
&\quad \left. - (\mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1}^l|\mathbf{x}_0^l, t) \| p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)) + \mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1}^l|\mathbf{x}_0^l, t) \| p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l))) \right)
\end{aligned}$$

Using the Gaussian parameterization of the reverse process (Eq. (1)), the above loss simplifies to:

$$L_1(\theta) \leq -\mathbb{E}_{t, \epsilon^w, \epsilon^l} \log \sigma \left(-\beta T \omega(\lambda_t) (\|\epsilon^w - \epsilon_\theta(\mathbf{x}_t^w, t)\|^2 - \|\epsilon^w - \epsilon_{\text{ref}}(\mathbf{x}_t^w, t)\|^2 - (\|\epsilon^l - \epsilon_\theta(\mathbf{x}_t^l, t)\|^2 - \|\epsilon^l - \epsilon_{\text{ref}}(\mathbf{x}_t^l, t)\|^2)) \right)$$

where $\epsilon^w, \epsilon^l \sim \mathcal{N}(0, I)$, $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$ thus $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon$. Same as Eq. (2), $\lambda_t = \alpha_t^2 / \sigma_t^2$ is a signal-to-noise ratio term [20], in practice, the reweighting assigns each term the same weight [17].

An alternative approximation Note that for Eq. (18) we utilize $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$ to approximate $p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0)$. For each step, it is to use $q(\mathbf{x}_{t-1,t}|\mathbf{x}_0)$ to approximate $p_\theta(\mathbf{x}_{t-1,t}|\mathbf{x}_0)$. Alternatively, we also propose to use $q(\mathbf{x}_t|\mathbf{x}_0)p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ for approximation. And this approximation yields lower error because $\mathbb{D}_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_0)p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \| p_\theta(\mathbf{x}_{t-1,t}|\mathbf{x}_0)) = \mathbb{D}_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_0) \| p_\theta(\mathbf{x}_t|\mathbf{x}_0)) < \mathbb{D}_{\text{KL}}(q(\mathbf{x}_{t-1,t}|\mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1,t}|\mathbf{x}_0))$.

$$\begin{aligned}
L_{\text{DPO-Diffusion}}(\theta) &= -\log \sigma \left(\beta \mathbb{E}_{\mathbf{x}_{1:T}^w \sim p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0^w), \mathbf{x}_{1:T}^l \sim p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0^l)} \left[\log \frac{p_\theta(\mathbf{x}_{0:T}^w)}{p_{\text{ref}}(\mathbf{x}_{0:T}^w)} - \log \frac{p_\theta(\mathbf{x}_{0:T}^l)}{p_{\text{ref}}(\mathbf{x}_{0:T}^l)} \right] \right) \\
&= -\log \sigma \left(\beta \mathbb{E}_{\mathbf{x}_{1:T}^w \sim p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0^w), \mathbf{x}_{1:T}^l \sim p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0^l)} \left[\sum_{t=1}^T \log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right) \\
&= -\log \sigma \left(\beta \mathbb{E}_{\mathbf{x}_{1:T}^w \sim p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0^w), \mathbf{x}_{1:T}^l \sim p_\theta(\mathbf{x}_{1:T}|\mathbf{x}_0^l)} T \mathbb{E}_t \left[\log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right) \\
&= -\log \sigma \left(\beta T \mathbb{E}_t \mathbb{E}_{\mathbf{x}_{t-1,t}^w \sim p_\theta(\mathbf{x}_{t-1,t}|\mathbf{x}_0^w), \mathbf{x}_{t-1,t}^l \sim p_\theta(\mathbf{x}_{t-1,t}|\mathbf{x}_0^l)} \left[\log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right).
\end{aligned}$$

By approximating $p_\theta(x_{t-1,t}|\mathbf{x}_0)$ with $q(\mathbf{x}_t|\mathbf{x}_0)p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$, we have

$$\begin{aligned} L_2(\theta) &= -\log \sigma \left(\beta T \mathbb{E}_t \mathbb{E}_{\mathbf{x}_{t-1,t}^w \sim q(\mathbf{x}_t|\mathbf{x}_0^w), \mathbf{x}_{t-1,t}^l \sim q(\mathbf{x}_t|\mathbf{x}_0^l)} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t^w), \mathbf{x}_{t-1,t}^l \sim q(\mathbf{x}_t|\mathbf{x}_0^l)} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t^l) \left[\log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right) \\ &= -\log \sigma \left(\beta T \mathbb{E}_{t, \mathbf{x}_t^w \sim q(\mathbf{x}_t|\mathbf{x}_0^w), \mathbf{x}_t^l \sim q(\mathbf{x}_t|\mathbf{x}_0^l)} \mathbb{E}_{\mathbf{x}_{t-1}^w \sim p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t^w), \mathbf{x}_{t-1}^l \sim p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t^l)} \left[\log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right). \end{aligned}$$

By Jensen's inequality, we have

$$\begin{aligned} L_2(\theta) &\leq -\mathbb{E}_{t, \mathbf{x}_t^w \sim q(\mathbf{x}_t|\mathbf{x}_0^w), \mathbf{x}_t^l \sim q(\mathbf{x}_t|\mathbf{x}_0^l)} \log \sigma \left(\beta T \mathbb{E}_{\mathbf{x}_{t-1}^w \sim p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t^w), \mathbf{x}_{t-1}^l \sim p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t^l)} \left[\log \frac{p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)}{p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)} - \log \frac{p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)}{p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)} \right] \right) \\ &= -\mathbb{E}_{t, \mathbf{x}_t^w \sim q(\mathbf{x}_t|\mathbf{x}_0^w), \mathbf{x}_t^l \sim q(\mathbf{x}_t|\mathbf{x}_0^l)} \log \sigma \left(\beta T \left(\mathbb{D}_{\text{KL}}(p_\theta(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w) \| p_{\text{ref}}(\mathbf{x}_{t-1}^w|\mathbf{x}_t^w)) - \mathbb{D}_{\text{KL}}(p_\theta(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l) \| p_{\text{ref}}(\mathbf{x}_{t-1}^l|\mathbf{x}_t^l)) \right) \right) \end{aligned}$$

Using the Gaussian parameterization of the reverse process (Eq. (1)), the above loss simplifies to:

$$L_2(\theta) = -\mathbb{E}_{t, \epsilon^w, \epsilon^l} \log \sigma \left(-\beta T \omega(\lambda_t) \left(\|\epsilon_\theta(\mathbf{x}_t^w, t) - \epsilon_{\text{ref}}(\mathbf{x}_t^w, t)\|^2 - \|\epsilon_\theta(\mathbf{x}_t^l, t) - \epsilon_{\text{ref}}(\mathbf{x}_t^l, t)\|^2 \right) \right)$$

where $\epsilon^w, \epsilon^l \sim \mathcal{N}(0, I)$, $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$ thus $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon$. Same as Eq. (2), $\lambda_t = \alpha_t^2 / \sigma_t^2$ is a signal-to-noise ratio term [20], in practice, the reweighting assigns each term the same weight [17].

S3. Alternate Derivation: Reinforcement Learning Perspective

We can also derive our objective as a multi-step RL approach, in the same setting as [6, 11]. A Markov Decision Process (MDP) is a tuple $(\mathcal{S}, \mathcal{A}, \rho_0, \mathcal{P}, \mathcal{R})$, where \mathcal{S} is the state space, \mathcal{A} is the action space, ρ_0 is an initial state distribution, \mathcal{P} is the transition dynamics and \mathcal{R} is the reward function. In this formulation, at each time step t a policy $\pi(a_t|s_t)$ observes a state $s_t \in \mathcal{S}$ and takes an action $a_t \in \mathcal{A}$. The environment then transitions to a next state $s_{t+1} \sim \mathcal{P}(s_{t+1}|s_t, a_t)$ and the returns a reward $\mathcal{R}(s_t, a_t)$. The goal of the policy is to maximize the total rewards it receives. Prior works [6, 11] map the denoising process in diffusion model generation to this formulation via:

$$\begin{aligned} \mathbf{s}_t &\triangleq (\mathbf{c}, \mathbf{x}_t, t) \\ \mathbf{a}_t &\triangleq \mathbf{x}_t \\ \mathcal{P}(s_{t+1}|\mathbf{s}_t, \mathbf{a}_t) &\triangleq (\delta_{\mathbf{c}}, \delta_{t-1}, \delta_{\mathbf{x}_{t-1}}) \\ \rho(\mathbf{s}_0) &\triangleq (p(\mathbf{c}), \delta_T, \mathcal{N}(\mathbf{0}, \mathbf{I})) \\ \mathcal{R}(\mathbf{s}_t, \mathbf{a}_t) &= \begin{cases} r(\mathbf{c}, \mathbf{x}_0) & \text{if } t = 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (19)$$

where \mathbf{c} is the prompt \mathbf{x}_t is the time-step t noisy image and δ_y is the Dirac delta function with unit density at y . That is in this formulation we consider the denoising model as a policy, with each denoising step a step in an MDP. The objective of the policy is to maximize the reward (alignment with human preference) of the final image. In the derivation below, we drop the time step t for brevity. In this formulation the generative model is a policy and the denoising process is a rollout in an MDP with a sparse reward received for the final generated image. Following [11] we optimize the following objective

$$\mathbb{E}_{\mathbf{c} \sim \mathcal{D}, p_\theta} \left[\sum_{t=T}^0 r(\mathbf{c}, \mathbf{x}_t) - \beta D_{\text{KL}}[p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c}) \| p_{\text{ref}}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})] \right] \quad (20)$$

While prior works [6, 11] use policy gradient approaches to optimize this objective, we're going to use off-policy methods. Following Control as Variational Inference [23], we have the following

$$Q^*((\mathbf{x}_t, \mathbf{c}), \mathbf{x}_{t-1}) = r(\mathbf{c}, \mathbf{x}_t) + V^*(\mathbf{x}_{t-1}, \mathbf{c}) \quad (21)$$

$$V^*(\mathbf{x}_{t-1}, \mathbf{c}) = \beta \log \mathbb{E}_{p_{\text{ref}}} [\exp Q^*((\mathbf{x}_t, \mathbf{c}), \mathbf{x}_{t-1}) / \beta] \quad (22)$$

$$p^*(\mathbf{x}_{t-1} | (\mathbf{x}_t, \mathbf{c})) = p_{\text{ref}}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c}) e^{(Q^*((\mathbf{x}_t, \mathbf{c}), \mathbf{x}_{t-1}) - V^*(\mathbf{x}_t, \mathbf{c})) / \beta} \quad (23)$$

where V^* is the optimal value function and Q^* is the optimal state-action value function (in our definition of the denoising MDP, the policy is stochastic, but the dynamics is deterministic). Also notice that in Eq. 23 the equation is exact since the right-hand side integrates to 1. We then consider the inverse soft Bellman operator [13] and have the following

$$r(\mathbf{c}, \mathbf{x}_t) = V^*(\mathbf{x}_{t-1}, \mathbf{c}) - Q^*((\mathbf{x}_t, \mathbf{c}), \mathbf{x}_{t-1}) \quad (24)$$

However, from Eq. 23 we have

$$Q^*((\mathbf{x}_t, \mathbf{c}), \mathbf{x}_{t-1}) - V^*(\mathbf{x}_t, \mathbf{c}) = \log \frac{p^*(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})}{p_{\text{ref}}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})} \quad (25)$$

substituting in Eq. 24 we obtain:

$$r(\mathbf{c}, \mathbf{x}_t) = V^*(\mathbf{x}_{t-1}, \mathbf{c}) + \log \frac{p^*(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})}{p_{\text{ref}}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})} - V^*(\mathbf{x}_t, \mathbf{c}) \quad (26)$$

Using a telescoping sum through the diffusion chain we are left with

$$r(\mathbf{c}, \mathbf{x}_0) = \sum_{t=0}^T \log \frac{p^*(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})}{p_{\text{ref}}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c})} - V^*(\mathbf{x}_T, \mathbf{c}) \quad (27)$$

since by definition all intermediate rewards are zero. If we assume both diffusion chains start from the same state and plug this result into the preference formulation of Eq. 3 we obtain the objective of Eq. 11. Here we optimize the same objective as prior works [6, 11], but instead of a policy gradient approach we derive our objective as an off-policy learning problem in the same MDP. This not only simplifies the algorithm significantly, but justifies our sampling choices in Eq. 13 and we do not have to sample through the entire diffusion chain.

S4. Alternative Derivation: Noise-Aware Preference Model

Paralleling the original DPO formulation we consider a policy trained on maximizing the likelihood of $p(x_0 | \mathbf{c}, t, x_{\text{obs}})$ where x_{obs} is a noised version of x_0 . Here x_0 is an image, \mathbf{c} is a text caption, t is a noising scale, and x_{obs} is a corruption (noised version) of x_0 . We initialize from a reference diffusion policy p_{ref} . We aim to optimize the same RL objective of Eq. (5), reprinted here for convenience:

$$\max_{p_\theta} \mathbb{E}_{\mathbf{c} \sim \mathcal{D}, \mathbf{x}_0 \sim p_\theta(\mathbf{x}_0 | \mathbf{c})} [r(\mathbf{c}, \mathbf{x}_0)] - \beta \mathbb{D}_{\text{KL}} [p_\theta(\mathbf{x}_0 | \mathbf{c}) \| p_{\text{ref}}(\mathbf{x}_0 | \mathbf{c})] \quad (28)$$

Our policy has additional conditioning (t, x_{obs}) . The latter is a noised version of x_0 . Define the space of noising operators at time t as Q_t where $q_t \sim Q_t$ with $q_t(x_0) = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon_{q_t}$, $\epsilon_{q_t} \sim N(0, I)$. Here q_t refers to the linear transform corresponding with a specific gaussian draw $\sim N(0, I)$ and the set of q_t is Q_t . In general at some time level t we have $y_{\text{obs}} = q_t(x_0)$ for some $q_t \sim Q_t$ so can write the conditioning as $p(x_0 | \mathbf{c}, t, q_t(y))$. We rewrite Eq. (28) as

$$\max_{p_\theta} \mathbb{E}_{\mathbf{c} \sim \mathcal{D}, \mathbf{x}_0 \sim p_\theta^{(\text{gen})}(\mathbf{x}_0 | \mathbf{c}), t \sim \mathcal{U}\{0, T\}, q_t \sim Q_T} (r_\phi(\mathbf{c}, x_0) - \beta \mathbb{D}_{\text{KL}} [p_\theta(x_0 | \mathbf{c}, t, q_t(x_0)) \| p_{\text{ref}}(x_0 | \mathbf{c}, t, q_t(x_0))]) \quad (29)$$

$p^{(\text{gen})}$ denoting the generative process associated with p as a diffusion model. Note that the reward model is the same formulation as in DPO. The optimal policy now becomes

$$p_\theta^*(x_0 | \mathbf{c}, t, q_t(x_0)) = \frac{1}{Z(\mathbf{c}, t, q_t)} p_{\text{ref}}(x_0 | \mathbf{c}, t, q_t(x_0)) \exp \left(\frac{1}{\beta} r(\mathbf{c}, x_0) \right) \quad (30)$$

with Z a partition over captions, timesteps, and noising draws. Rearranging for $r(\mathbf{c}, x_0)$ now yields

$$r(\mathbf{c}, x_0) = \beta \log \frac{p_\theta^*(x_0 | \mathbf{c}, t, q_t)}{p_{\text{ref}}(x_0 | \mathbf{c}, t, q_t)} + \beta \log Z(\mathbf{c}, t, q_t), \quad \forall t, q_t \quad (31)$$

We have not changed the reward model formulation at all, but our policies have extra conditioning as input (which ideally the likelihoods are constant with respect to). Putting this formulation into the original Bradley-Terry framework of Eq. (3) (re-printed here)

$$p_{\text{BT}}(\mathbf{x}_0^w \succ \mathbf{x}_0^l | \mathbf{c}) = \sigma(r(\mathbf{c}, \mathbf{x}_0^w) - r(\mathbf{c}, \mathbf{x}_0^l)) \quad (32)$$

results in the objective:

$$\mathcal{L}_{\text{DPO}}(p_\theta; p_{\text{ref}}) = -\mathbb{E}_{(\mathbf{x}_0^w, \mathbf{x}_0^l \sim p^{(\text{gen})}(\mathbf{c}); \mathbf{c} \sim \mathcal{D}, t \sim \mathcal{U}\{0, T\}; q_t \sim Q_t)} \left[\log \sigma \left(\beta \log \frac{p_\theta(\mathbf{x}_0^w | \mathbf{c}, t, q_t(\mathbf{x}_0^w))}{p_{\text{ref}}(\mathbf{x}_0^w | \mathbf{c}, t, q_t(\mathbf{x}_0^w))} - \beta \log \frac{p_\theta(\mathbf{x}_0^l | \mathbf{c}, t, q_t(\mathbf{x}_0^l))}{p_{\text{ref}}(\mathbf{x}_0^l | \mathbf{c}, t, q_t(\mathbf{x}_0^l))} \right) \right] \quad (33)$$

We now consider how to compute these likelihoods. Using the notation $a_t = \sqrt{\alpha_t}$ and $b_t = \sqrt{1 - \alpha_t}$ as shorthand for commonly-used diffusion constants (α are defined as in DDIM[43]) we have

$$x_{\text{obs}} = q_t(x_0) = a_t x_0 + b_t \epsilon, \epsilon \sim \mathcal{N}(0, I) \quad (34)$$

We use Eq. 57 from DDIM[43] (along with their definition of σ_t):

$$p(x_0 | x_t) = \mathcal{N}(x_0^{\text{pred}}, \sigma_t^2 I) \quad (35)$$

Our x_0^{pred} is:

$$x_0^{\text{pred}} = \frac{x_{\text{obs}} - b_t \epsilon_\theta^{\text{pred}}}{a_t} = \frac{a_t x_0 + b_t \epsilon - b_t \epsilon_\theta^{\text{pred}}}{a_t} = x_0 + \frac{b_t}{a_t} (\epsilon - \epsilon_\theta^{\text{pred}}) \quad (36)$$

Here $\epsilon_\theta^{\text{pred}}$ is the output of $\epsilon_\theta(\mathbf{c}, t, x_{\text{obs}})$ Making the conditional likelihood:

$$p_\theta(x_0 | \mathbf{c}, t, x_{\text{obs}}) = \mathcal{N}(x_0; x_0 + \frac{b_t}{a_t} (\epsilon - \epsilon_\theta^{\text{pred}}), \sigma_t^2 I) = \frac{1}{(2\pi\sigma_t^2)^{d/2}} e^{-\frac{b_t^2}{2a_t^2\sigma_t^2} \|\epsilon - \epsilon_\theta^{\text{pred}}\|_2^2} \quad (37)$$

For convenience we define

$$z_t = \frac{1}{(2\pi\sigma_t^2)^{d/2}} \quad (38)$$

$$SE = \|\epsilon - \epsilon_{\text{pred}}\|_2^2 \quad (39)$$

We will decorate the latter quantity (SE) with sub/superscripts later. For now we get:

$$p_\theta(x_0 | \mathbf{c}, t, x_{\text{obs}}) = z_t e^{-\frac{b_t^2}{2a_t^2\sigma_t^2} SE} \quad (40)$$

We see to minimize

$$\mathbb{E}_{(\mathbf{x}_0^w, \mathbf{x}_0^l \sim p^{(\text{gen})}(\mathbf{c}); \mathbf{c} \sim \mathcal{D}, t \sim \mathcal{U}\{0, T\}; q_t \sim Q_t)} -\log \sigma \left(\beta \left(\log \frac{p_\theta(\mathbf{x}_0^w | \mathbf{c}, t, q_t(\mathbf{x}_0^w))}{p_{\text{ref}}(\mathbf{x}_0^w | \mathbf{c}, t, q_t(\mathbf{x}_0^w))} - \log \frac{p_\theta(\mathbf{x}_0^l | \mathbf{c}, t, q_t(\mathbf{x}_0^l))}{p_{\text{ref}}(\mathbf{x}_0^l | \mathbf{c}, t, q_t(\mathbf{x}_0^l))} \right) \right) = \quad (41)$$

$$\mathbb{E}_{(\mathbf{x}_0^w, \mathbf{x}_0^l \sim p^{(\text{gen})}(\mathbf{c}); \mathbf{c} \sim \mathcal{D}, t \sim \mathcal{U}\{0, T\}; q_t \sim Q_t)} -\log \sigma \left(\beta \left(\log \frac{z_t e^{-\frac{b_t^2}{2a_t^2\sigma_t^2} SE_\theta^{(w)}}}{z_t e^{-\frac{b_t^2}{2a_t^2\sigma_t^2} SE_{\text{ref}}^{(w)}}} - \log \frac{z_t e^{-\frac{b_t^2}{2a_t^2\sigma_t^2} SE_\theta^{(l)}}}{z_t e^{-\frac{b_t^2}{2a_t^2\sigma_t^2} SE_{\text{ref}}^{(l)}}} \right) \right) \quad (42)$$

Here we use $SE_\psi^{(d)} = \|\epsilon_{q_t} - \psi(\mathbf{c}, t, q_t(x_0^d))\|_2^2$ to denote the L2 error in the noise prediction of model ψ operating on the noisy $q_t(x_0^d)$ with corresponding conditioning (\mathbf{c}, t) ($d \in \{w, l\}$). Here the model associated with SE_{ref}^* is the model of the reference policy p_{ref} . Note that these SE terms are the standard diffusion training objective from Eq. (2). Continuing to simplify the above yields:

$$-\log \sigma \left(\beta \left(\log \frac{z_t e^{-\frac{b_t^2}{2a_t^2 \sigma_t^2} SE_\theta^{(w)}}}{z_t e^{-\frac{b_t^2}{2a_t^2 \sigma_t^2} SE_{\text{ref}}^{(w)}}} - \log \frac{z_t e^{-\frac{b_t^2}{2a_t^2 \sigma_t^2} SE_\theta^{(l)}}}{z_t e^{-\frac{b_t^2}{2a_t^2 \sigma_t^2} SE_{\text{ref}}^{(l)}}} \right) \right) \quad (43)$$

$$= -\log \sigma \left(-\beta \frac{b_t^2}{2a_t^2 \sigma_t^2} \left((SE_\theta^{(w)} - SE_{\text{ref}}^{(w)}) - (SE_\theta^{(l)} - SE_{\text{ref}}^{(l)}) \right) \right) \quad (44)$$

We can simplify the coefficient:

$$\frac{b_t^2}{a_t^2 \sigma_t^2} = \frac{1 - \alpha_t}{\alpha_t} \frac{1}{\sigma_t^2} = \frac{\sigma_{t+1}^2}{\sigma_t^2} \approx 1 \quad (45)$$

Resulting in objective

$$\approx_{x, y_w, y_l \sim D; t; \epsilon \sim \mathcal{N}(0, I)} \mathbb{E} -\log \sigma \left(-\frac{\beta}{2} \left((SE_\theta^{(w)} - SE_{\text{ref}}^{(w)}) - (SE_\theta^{(l)} - SE_{\text{ref}}^{(l)}) \right) \right) \quad (46)$$

Up to the approximation of Eq. (45) this is the equivalent to Eq. (33). The log of the likelihood ratios simply take on the elegant form of a difference in diffusion training losses. Due to the equation negatives and log σ being a monotonic increasing function, by minimizing Eq. (46) we are aiming to minimize the inside term

$$\left((SE_\theta^{(w)} - SE_{\text{ref}}^{(w)}) - (SE_\theta^{(l)} - SE_{\text{ref}}^{(l)}) \right) \quad (47)$$

This can be done by minimizing $SE_\theta^{(w)}$ or maximizing $SE_\theta^{(l)}$, with the precise loss value depending on how these compare to the reference errors $SE_{\text{ref}}^{(w)}, SE_{\text{ref}}^{(l)}$. The asymmetry of the log σ function allows β to control the penalty for deviating from the reference distribution. A high β results in a highly asymmetric distribution, disproportionately penalizing low $SE_\theta^{(l)}$ and high $SE_\theta^{(w)}$ and encouraging a p_θ to make less mistakes in implicitly scoring y_w, y_l by deviating less from the reference policy p_{ref} . We visualize the log σ curves in Figure S1 for several values of β .

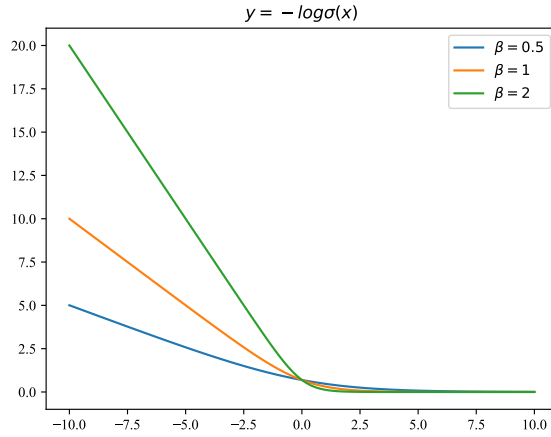


Figure S1. Visualization of $y = -\log \sigma(\beta x)$

S4.1. Reward Estimation

Finally, we note that in this formulation that if we wish to compute the noise-aware reward difference $r(\mathbf{c}, x_0^A) - r(\mathbf{c}, x_0^B)$, referring to Eq. (31) this now has form

$$r(\mathbf{c}, x_0^A) - r(\mathbf{c}, x_0^B) = [\beta(SE_\theta^A - SE_{\text{ref}}^A) + \beta \log Z(\mathbf{c}, t, q_t)] - [\beta(SE_\theta^B - SE_{\text{ref}}^B) + \beta \log Z(\mathbf{c}, t, q_t)], \forall \mathbf{c}, t, q_t \quad (48)$$

$$= \beta[(SE_\theta^A - SE_{\text{ref}}^A) - (SE_\theta^B - SE_{\text{ref}}^B)], \forall \mathbf{c}, t, q_t \quad (49)$$

$$(50)$$

Which means for two images (x_0^A, x_0^B) with the same conditioning \mathbf{c} we can estimate the reward difference using Eq. (48). When doing this it improves the estimate to average over multiple draws (t, q_t) . We use this method in Table S2.

S5. β Ablation



Figure S2. Median PickScores for generations on the Pick-a-Pic v2 validation set for different choices of β

For β far below the displayed regime, the diffusion model degenerates into a pure reward scoring model. Much greater, and the KL-divergence penalty greatly restricts any appreciable adaptation.

S6. Further SFT Discussions

We also partially attribute this difference in effectiveness of SFT to the gap in pretraining vs. downstream task considered in the original DPO paper [33] vs. our work. On two of the DPO LLM tasks (**sentiment generation**, **single-turn dialogue**), generic off-the-shelf autoregressive language models are tuned on specific tasks in the SFT stage. In the final setting, **summarization**, the SFT model has been pretrained on a similar task/dataset. In this case, finetuning on the “preferred” dataset (preferred-FT) baseline performs comparably to the SFT initialization.

This final setting is most analogous to that of Diffusion-DPO. The generic pretraining, task, and evaluation setting are all text-to-image generation. There is no task-specific domain gap and all of the settings are open-vocabulary with a broad range of styles. As such, our findings are similar to that of **summarization** in [33] where an already task-tuned model does not benefit from preferred finetuning.

S7. Additional Automated Metrics

Automated metrics on Pick-a-Pic validation captions are shown in Figure S3 for DPO-SDXL. The y-axis measures the fraction of head-to-head generation comparisons for a prompt that DPO-SDXL scores higher than the baseline SDXL.

S8. PickScore Rejection Sampling

Rejection sampling was used in [21] as a powerful inference-time tool. 100 samples were drawn from variants of a prompt and PickScore-ranked, with the highest scored images being compared to a single random draw. PickScore selections were human-preferred 71.4% of the time. We compare using additional compute at inference vs. additional training in Figure S4. We plot the expected PickScore win rate of n draws from the reference model against a single draw from the learned (DPO) model. The mean inference compute for baseline rejection sampling to surpass the DPO-trained model is $10\times$ higher in both cases. For 7% (SDXL) and 16% (SD1.5) of the prompts even 100 draws is insufficient.

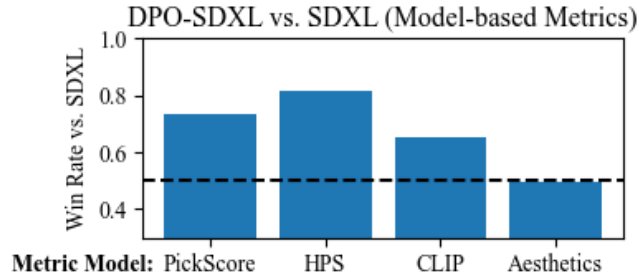


Figure S3. Box plots of automated metrics vs. SDXL baseline. All 500 unique prompts from PickScore validation set.

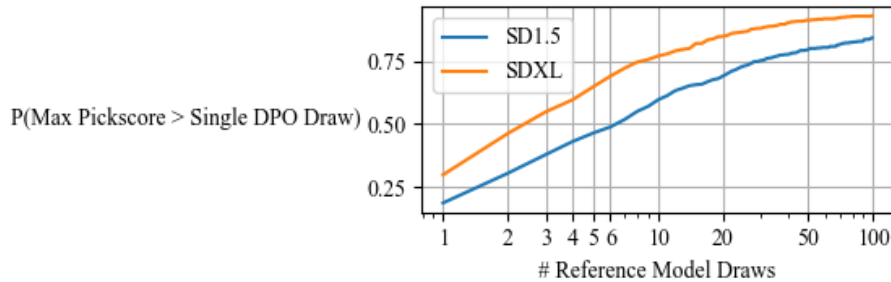


Figure S4. The number of draws from the reference model vs. the probability that maximum PickScore of the draws exceeds a single DPO generation. 500 PickScore validation prompts used. Mean (including 100s)/Median: SDXL (13.7, 3), SD1.5 (25.6, 7).

S9. Pseudocode for Training Objective

```
def loss(model, ref_model, x_w, x_l, c, beta):
    """
    # This is an example psuedo-code snippet for calculating the Diffusion-DPO loss
    # on a single image pair with corresponding caption

    model: Diffusion model that accepts prompt conditioning c and time conditioning t
    ref_model: Frozen initialization of model
    x_w: Preferred Image (latents in this work)
    x_l: Non-Preferred Image (latents in this work)
    c: Conditioning (text in this work)
    beta: Regularization Parameter

    returns: DPO loss value
    """
    timestep = torch.randint(0, 1000)
    noise = torch.randn_like(x_w)
    noisy_x_w = add_noise(x_w, noise, t)
    noisy_x_l = add_noise(x_l, noise, t)

    model_w_pred = model(noisy_x_w, c, t)
    model_l_pred = model(noisy_x_l, c, t)
    ref_w_pred = ref(noisy_x_w, c, t)
    ref_l_pred = ref(noisy_x_l, c, t)
```

```
model_w_err = (model_w_pred - noise).norm().pow(2)
model_l_err = (model_l_pred - noise).norm().pow(2)
ref_w_err = (ref_w_pred - noise).norm().pow(2)
ref_l_err = (ref_l_pred - noise).norm().pow(2)

w_diff = model_w_err - ref_w_err
l_diff = model_l_err - ref_l_err

inside_term = -1 * beta * (w_diff - l_diff)

loss = -1 * log(sigmoid(inside_term))

return loss
```

S10. Additional Qualitative Results

In Figure S6 we present generations from DPO-SDXL on complex prompts from DALLE3 [15]. Other generations for miscellaneous prompts are shown in Figure S5. In Fig. S7 and 8 we display qualitative comparison results from HPSv2 with random seeds from our human evaluation for prompt indices 200, 600, 1000, 1400, 1800, 2200, 2600, 3000.



Figure S5. DPO-SDXL gens on miscellaneous prompts Prompts (clockwise) (1) A bulldog mob boss, moody cinematic feel (2) A old historical notebook detailing the discovery of unicorns (3) A purple raven flying over a forest of fall colors, imaginary documentary (4) Small dinosaurs shopping in a grocery store, oil painting (5) A wolf wearing a sheep halloween costume going trick-or-treating at the farm (6) A mummy studying hard in the library for finals, head in hands



Figure S6. DPO-SDXL gens on prompts from DALLE3 [15] Prompts: (1): A swirling, multicolored portal emerges from the depths of an ocean of coffee, with waves of the rich liquid gently rippling outward. The portal engulfs a coffee cup, which serves as a gateway to a fantastical dimension. The surrounding digital art landscape reflects the colors of the portal, creating an alluring scene of endless possibilities. (2): In a fantastical setting, a highly detailed furry humanoid skunk with piercing eyes confidently poses in a medium shot, wearing an animal hide jacket. The artist has masterfully rendered the character in digital art, capturing the intricate details of fur and clothing texture"



Figure S7. Prompts: (1) A kangaroo wearing an orange hoodie and blue sunglasses stands on the grass in front of the Sydney Opera House, holding a sign that says Welcome Friends. (2) Anime Costa Blanca by Studio Ghibli. (3) There is a secret museum of magical items inside a crystal greenhouse palace filled with intricate bookshelves, plants, and Victorian style decor. (4) A depiction of Hermione Granger from the Harry Potter series as a zombie.

SDXL

DPO-SDXL



Figure S8. (1) A portrait art of a necromancer, referencing DND and War craft. (2) Monalisa painting a portrait of Leonardo Da Vinci. (3) There is a cyclist riding above all the pigeons. (4) A woman holding two rainbow slices of cake.