

# A TAILORED FRAMEWORK FOR ALIGNING DIFFUSION MODELS WITH HUMAN PREFERENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The direct preference optimization (DPO) method has shown success in aligning text-to-image diffusion models with human preference. Previous approaches typically assume a consistent preference label between final generated images and their corresponding noisy samples at intermediate steps, and directly apply DPO to these noisy samples for fine-tuning. However, we identify a significant issue with this consistency assumption, as directly applying DPO to noisy samples from different generation trajectories based on final preference order may disrupt the optimization process. We first demonstrate the issues inherent in previous methods from two perspectives: *gradient direction* and *preference order*, and then propose a **Tailored Preference Optimization** (TailorPO) framework for aligning diffusion models with human preference, underpinned by some theoretical insights. Our approach directly ranks the preference order of intermediate noisy samples based on their step-wise reward, and effectively resolves the optimization direction issues through a simple yet efficient design. Additionally, to the best of our knowledge, we are the first to consider the distinct structure of diffusion models and leverage the gradient guidance in preference aligning to enhance the optimization effectiveness. Experimental results demonstrate that our method significantly improves the model’s ability to generate aesthetically pleasing and human-preferred images.

## 1 INTRODUCTION

Direct preference optimization (DPO), which fine-tunes the model on paired data to align the model generations with human preferences, has demonstrated its success in large language models (LLMs) (Rafailov et al., 2023). Recently, researchers generalized this method to diffusion models for text-to-image generation (Black et al., 2024; Yang et al., 2024a; Wallace et al., 2024). Given a pair of images generated from the same prompt and a ranking of human preference for them, DPO aims to increase the probability of generating the preferred sample while decreasing the probability of generating another sample, which enables the model to generate more visually appealing and aesthetically pleasing images that better align with human preferences.

Specifically, previous researchers (Yang et al., 2024a) leverage the *trajectory-level* preference to rank the generated samples. As shown in Figure 1(a), given a text prompt  $c$ , they first sample a pair of denoising trajectories  $[x_T^0, \dots, x_0^0]$  and  $[x_T^1, \dots, x_0^1]$  from the diffusion model, and then rank them according to the human preference on the final generated images  $x_0^0$  and  $x_0^1$ . It is assumed that the *preference order of  $(x_0^0, x_0^1)$ , at the end of the generation trajectory, can consistently represent the preference order of  $(x_t^0, x_t^1)$  at all intermediate steps  $t$* . Then, the DPO loss function is implemented using the generation probabilities  $p(x_{t-1}^0|x_t^0, c)$  and  $p(x_{t-1}^1|x_t^1, c)$  at each step  $t$  to fine-tune the diffusion model, which is called the *step-level* optimization.

However, we notice that the above trajectory-level preference ranking and the step-level optimization are not fully compatible in diffusion models. **First**, the trajectory-level preference ranking (*i.e.*, the preference order of final outputs  $(x_0^0, x_0^1)$  of trajectories) does not accurately reflect the preference order of  $(x_t^0, x_t^1)$  at intermediate steps. Considering the inherent randomness in the denoising process, simply assigning the preference of final outputs to all the intermediate steps will detrimentally affect the preference optimization performance. **Second**, the generation probabilities  $p(x_{t-1}^0|x_t^0, c)$  and  $p(x_{t-1}^1|x_t^1, c)$  in two different trajectories are conditioned on different inputs, and this causes the optimization direction to be significantly affected by the difference between the inputs. In particular,

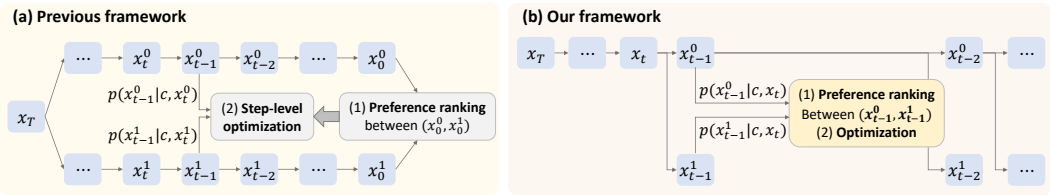


Figure 1: Framework overview of (a) previous method and (b) TailorPO. In the previous method, the preference order is determined based on final outputs and used to guide the optimization of intermediate noisy samples in different generation trajectories. In contrast, we generate noisy samples from the same input  $x_t$  and directly rank their preference order for optimization.

if  $x_t^0$  and  $x_t^1$  are located in the same linear subspace of the diffusion model, then the optimization of DPO probably increases the output probability of the dis-preferred samples. We conducted a detailed theoretical analysis of these issues in Section 3.2.

Therefore, in this paper, we propose a **Tailored Preference Optimization (TailorPO)** framework to fine-tune diffusion models with DPO, which addresses the aforementioned challenges. As Figure 1(b) shows, we generate two different noisy samples  $(x_{t-1}^0, x_{t-1}^1)$  from the same input  $x_t$  at each denoising step. Then, we directly rank the preference order of two samples  $(x_{t-1}^0, x_{t-1}^1)$  based on their step-wise reward. To this end, one of the most straightforward methods is to directly evaluate the reward of these noisy samples using a reward model. However, most existing reward models are trained on natural images and are not applicable to noisy samples. To address this issue, we formulate the denoising process as a Markov decision process (MDP) and derive a simple yet effective measurement for the preference reward of noisy samples. Then, given the preference order, we utilize  $p(x_{t-1}^0|x_t, c)$  and  $p(x_{t-1}^1|x_t, c)$  to compute the DPO loss function for fine-tuning. In this way, the gradient direction is proven to increase the probability of generating preferred samples while decreasing the probability of generating dis-preferred samples.

Moreover, we notice that TailorPO generates paired samples from the same  $x_t$ , potentially causing two samples to be similar in late denoising steps with large  $t$ . Such similarity may reduce the diversity of paired samples, thereby impacting the effectiveness of the DPO-based method. To mitigate this issue, we propose to enhance the diversity of noisy samples by increasing their reward gap. Specifically, we employ gradient guidance (Guo et al., 2024) to generate paired samples, leveraging the gradient of differentiable reward models to increase the reward of preferred noisy samples. This strategy, termed *TailorPO-G*, further improves the effectiveness of our TailorPO framework.

In experiments, we fine-tune Stable Diffusion v1.5 using TailorPO and TailorPO-G to enhance its ability to generate images that achieve elevated aesthetic scores and align with human preference. Additionally, we evaluate TailorPO on user-specific preferences, such as image compressibility. The experimental results indicate that diffusion models fine-tuned with TailorPO and TailorPO-G yield higher reward scores compared to those fine-tuned with other RLHF and DPO-style methods.

**Contributions** of this paper can be summarized as follows. (1) Through theoretical analysis and experimental validation, we demonstrate the mismatch between the trajectory-level ranking and the step-level optimization in existing DPO methods for diffusion models. (2) Based on these insights, we propose TailorPO, a simple DPO framework tailored to the unique denoising structure of diffusion models. To the best of our knowledge, this is the first framework that explicitly considers the properties of diffusion models for DPO. Experimental results have demonstrated that TailorPO significantly improves the model’s ability to generate human-preferred images. (3) Furthermore, inspired by the success of gradient guidance in adapting model outputs towards user-specified objectives, we incorporate gradient guidance of differentiable reward models in TailorPO-G to increase the diversity of training samples for fine-tuning to further enhance performance.

## 2 RELATED WORKS

**Diffusion models.** As a new class of generative models, diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021) transform Gaussian noises into images (Dhariwal &

Nichol, 2021; Ho et al., 2022b; Nichol et al., 2022; Rombach et al., 2022), audios (Liu et al., 2023), videos (Ho et al., 2022a; Singer et al., 2023), 3D shapes (Zeng et al., 2022; Poole et al., 2023; Gu et al., 2023), and robotic trajectories (Janner et al., 2022; Chen et al., 2024) through an iterative denoising process. Dhariwal & Nichol (2021) and Ho & Salimans (2022) further propose the classifier guidance and classifier-free guidance respectively to align the generated images with specific text descriptions for text-to-image synthesis.

**Learning diffusion models from human feedback.** Inspired by the success of reinforcement learning from human feedback (RLHF) in large language models (Ouyang et al., 2022; Bai et al., 2022; OpenAI, 2023), many reward models have been developed for images preference, including aesthetic predictor (Schuhmann et al., 2022), ImageReward (Xu et al., 2023), PickScore model (Kirstain et al., 2023), and HPSv2 (Wu et al., 2023). Based on these reward models, Lee et al. (2023), DPOK (Fan et al., 2023) and DDPO (Black et al., 2024) formulated the denoising process of diffusion models as a Markov decision process (MDP) and fine-tuned diffusion models using the policy-gradient method. DRaFT (Clark et al., 2024), and AlignProp (Prabhudesai et al., 2023) directly back-propagated the gradient of reward models through the sampling process of diffusion models for fine-tuning. In comparison, D3PO Yang et al. (2024a) and Diffusion DPO (Wallace et al., 2024) adapted the direct preference optimization (DPO) (Rafailov et al., 2023) on diffusion models and optimized model parameters at each denoising step. Considering the sequential nature of the denoising process, DenseReward (Yang et al., 2024b) assigned a larger weights for initial steps than later steps when using DPO. Most close to our work, SPO (Liang et al., 2024) also pointed out the problematic assumption about the preference consistency of intermediate noisy samples and final output images. However, they addressed this by training a step-wise reward model on another uncertain assumption. In comparison, we conduct a detailed analysis of the assumption and develop a new framework to improve the performance of DPO.

### 3 METHOD

#### 3.1 PRELIMINARIES

**Diffusion models.** Diffusion models contain a forward process and a reverse denoising process. In the forward process, given an input  $x_0$  sampled from the real distribution  $p_{\text{data}}$ , diffusion models gradually add Gaussian noises to  $x_0$  at each step  $t \in [1, T]$ , as follows:

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1} = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon \quad (1)$$

where  $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  denotes the Gaussian noise at step  $t$ .  $\alpha_{1:T}$  denotes the variance schedule and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ .

In the reverse denoising process, the diffusion model is trained to learn  $p(x_{t-1}|x_t)$  at each step  $t$ . Specifically, following (Song et al., 2021), the denoising step at step  $t$  is formulated as

$$x_{t-1} = \underbrace{\sqrt{\bar{\alpha}_{t-1}} \left( \frac{x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_{\theta}(x_t, t)}{\sqrt{\bar{\alpha}_t}} \right)}_{\hat{x}_0(x_t), \text{ predicted } x_0} + \underbrace{\sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2}\epsilon_{\theta}(x_t, t)}_{\text{direction pointing to } x_t} + \underbrace{\sigma_t\epsilon'_t}_{\text{random noise}} \quad (2)$$

where  $\epsilon_{\theta}(\cdot)$  is a noise prediction network with trainable parameters  $\theta$ , which aims to use  $\epsilon_{\theta}(x_t, t)$  to predict the noise  $\epsilon$  in Eq. (1) at each step  $t$ .  $\epsilon'_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is sampled from the standard Gaussian distribution. In fact,  $x_{t-1}$  is sampled from the estimated distribution  $\mathcal{N}(\mu_{\theta}(x_t), \sigma_t^2\mathbf{I})$ . According to the reverse process,  $\hat{x}_0(x_t) = (x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_{\theta}(x_t, t))/\sqrt{\bar{\alpha}_t}$  represents the predicted  $x_0$  at step  $x$ .

**Direct preference optimization (DPO) (Rafailov et al., 2023).** The DPO method is first proposed to fine-tune large language models to align with human preferences. Given a prompt  $x$ , two responses  $y_0$  and  $y_1$  are sampling from the generative model  $\pi_{\theta}$ , i.e.,  $y_0, y_1 \sim \pi_{\theta}(y|x)$ . Then,  $y_0$  and  $y_1$  are ranked based on human preferences or the outputs  $r(x, y_0)$  and  $r(x, y_1)$  of a pre-trained reward model  $r(\cdot)$ . Let  $y_w$  denote the preferred response in  $(y_0, y_1)$  and  $y_l$  denote the dis-preferred response. DPO optimizes parameters  $\theta$  in  $\pi_{\theta}$  by minimizing the following loss function.

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(x, y_w, y_l)} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right] \quad (3)$$

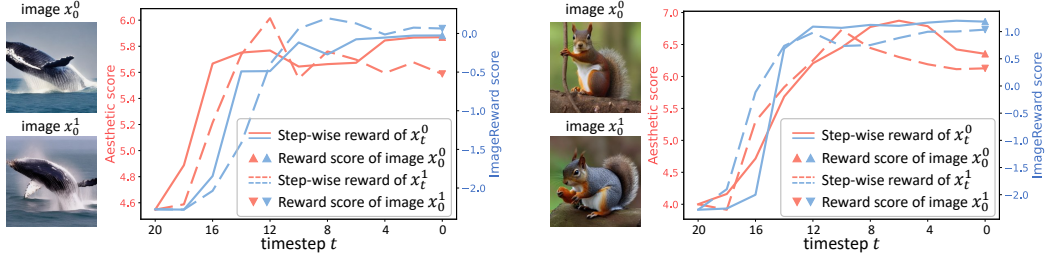


Figure 2: The preference order of intermediate noisy samples is not always consistent with the preference order of final output images, from both perspectives of the aesthetic score (red) and ImageReward score (blue).

where  $\sigma$  is the sigmoid function, and  $\beta$  is a hyper-parameter.  $\pi_{\text{ref}}$  represents the reference model, usually set as the pre-trained models before fine-tuning. The gradient of the above loss function on each pair of  $(x, y_w, y_l)$  with respect to the parameters  $\theta$  is as follows (Rafailov et al., 2023).

$$\nabla_{\theta} \mathcal{L}_{\text{DPO}}(\theta, x, y_w, y_l) = -f(x, y_w, y_l) (\nabla_{\theta} \log \pi_{\theta}(y_w|x) - \nabla_{\theta} \log \pi_{\theta}(y_l|x)) \quad (4)$$

where  $f(x, y_w, y_l) \triangleq \beta(1 - \sigma(\beta \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)}))$ . Therefore, the gradient of the DPO loss function increases the likelihood of the preferred response  $y_w$  and decreases the likelihood of the dis-preferred response  $y_l$ .

### 3.2 MISMATCH BETWEEN TRAJECTORY-LEVEL RANKING AND STEP-LEVEL OPTIMIZATION

In this section, we first revisit how existing works implement DPO for diffusion models, using D3PO (Yang et al., 2024a) as an example for explanation. Then, we identify the mismatch between their trajectory-level ranking and step-level optimization from two perspectives.

For a text-to-image diffusion model  $\pi_{\theta}$  parameterized by  $\theta$ , given a text prompt  $c$ , D3PO first samples a pair of generation trajectories  $[x_T^0, \dots, x_0^0]$  and  $[x_T^1, \dots, x_0^1]$ . Then, they compare the reward scores  $r(c, x_0^0)$  and  $r(c, x_0^1)$  of generated images, using the reward model  $r(\cdot)$ , and rank their preference order. The preferred image is denoted by  $x_0^w$  and the dis-preferred image is denoted by  $x_0^l$ . Then, as Figure 1(a) shows, it is assumed that the preference order of final images  $(x_0^0, x_0^1)$  represents the preference order of  $(x_t^0, x_t^1)$  at all intermediate steps  $t$ . Subsequently, the diffusion model is fine-tuned by minimizing the following DPO-like loss function at the step level.

$$\mathcal{L}_{\text{D3PO}}(\theta) = -\mathbb{E}_{(c, x_t^w, x_t^l, x_{t-1}^w, x_{t-1}^l)} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(x_{t-1}^w | x_t^w, c)}{\pi_{\text{ref}}(x_{t-1}^w | x_t^w, c)} - \beta \log \frac{\pi_{\theta}(x_{t-1}^l | x_t^l, c)}{\pi_{\text{ref}}(x_{t-1}^l | x_t^l, c)} \right) \right] \quad (5)$$

We argue that there are two critical issues in the aforementioned process and loss function, which we will elaborate on and prove through the theoretical analysis in the following sections.

**Inaccurate preference order.** The first obvious issue is that the preference order of final images  $x_0$  at the end of the trajectory cannot accurately reflect the preference order of noisy samples  $x_t$  at intermediate steps. Liang et al. (2024) demonstrated that early steps in the denoising process tend to handle layout, while later steps focus more on detailed textures. However, the preference order based on final images primarily reflects layout and composition preferences, misaligning with the function of later steps. Taking a step further, we rethink this problem from another perspective and formulate the reward at intermediate steps based on theoretical analysis.

Similar to (Yang et al., 2024a), we formulate the denoising process in a diffusion model as a Markov decision process (MDP), as follows.

$$\begin{aligned} S_t &\triangleq (c, x_{T-t}), A_t \triangleq x_{T-t-1}, R_t = R(S_t, A_t) \triangleq R((c, x_{T-t}), x_{T-t-1}) \\ P(S_{t+1}|S_t, A_t) &\triangleq (\delta_c, \delta_{x_{T-t-1}}), \pi(A_t|S_t) \triangleq \pi_{\theta}(x_{T-t-1}|x_{T-t}, c) \end{aligned} \quad (6)$$

where  $S_t, A_t, R_t, P(S_{t+1}|S_t, A_t)$ , and  $\pi(A_t|S_t)$  denote the state, action, reward, state transition probability, and the policy in MDP, respectively. In this finite MDP, the cumulative return at time  $t$

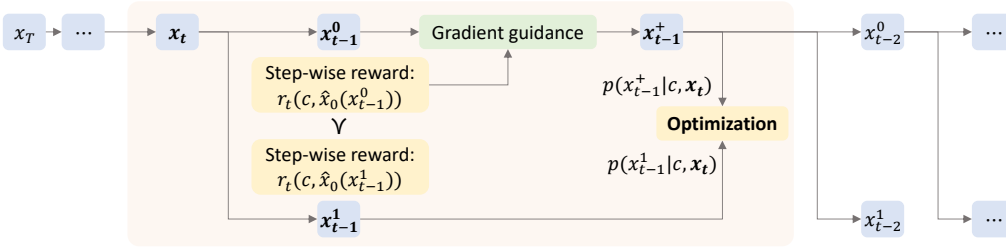


Figure 3: Framework of TailorPO. At each step  $t$ , we start from the same  $x_t$  to generate two noisy samples  $x_{t-1}^0$  and  $x_{t-1}^1$ . Subsequently, we compare their step-wise reward to determine their preference order. For the preferred sample, if the reward model is differentiable, we employ the gradient guidance to further increase its reward to obtain  $x_{t-1}^+$ . Then, we optimize the generating probability of preferred and dis-preferred samples. After the optimization at step  $t$ , the preferred sample is taken as the input  $x_{t-1}$  of the next step for later sampling and optimization.

can be defined as  $G_t = \sum_{k=t+1}^T R_k$ , and the action value function at time  $t$  is  $Q(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a]$ . For the denoising process of diffusion models, we simplify the cumulative return to the reward of the generated image, *i.e.*,  $G_t = R_T = r(c, x_0)$ . In this way, the action value function is simplified as follows.

$$Q(s, a) = \mathbb{E}[r(c, x_0) | S_t = (c, x_{T-t}), A_t = x_{T-t-1}] = \mathbb{E}[r(c, x_0) | c, x_{T-t-1}] \quad (7)$$

In other words, the quality of noisy samples  $x_{T-t-1}$  can be determined by the expected reward value of images generated by different trajectories starting from  $x_{T-t-1}$ . In contrast, the reward value  $r(c, x_0)$  of an image from a single trajectory does not represent the quality of the intermediate denoising action. Based on this analysis, we demonstrate that *the preference order of final images cannot accurately represent the preference order of intermediate noisy samples*.

To better illustrate this issue, we first propose a method for evaluating the quality of intermediate noisy samples, followed by an experimental validation using this method. The results shown in Figure 2 demonstrate that the preference order between a pair of intermediate samples  $x_t$  can sometimes conflict with the preference order between the corresponding denoised images  $x_0$ . This finding likewise provides evidence against the validity of the assumption employed in previous methods. The proposed evaluation method and our framework will be elaborated in the subsequent sections.

**Disturbed gradient direction.** Moreover, even if we obtain an accurate preference order of noisy samples at intermediate steps, the loss function in Eq. (5) still has limitations from the gradient perspective. To gain a mechanistic understanding of the above loss function, we compute its gradient with respect to parameters  $\theta$  as follows (please refer to Appendix A for the proof).

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\text{D3PO}}(\theta) &= -\mathbb{E} \left[ (f_t / \sigma_t^2) \cdot [(x_{t-1}^w - \mu_{\theta}(x_t^w))^T \nabla_{\theta} \mu_{\theta}(x_t^w) - (x_{t-1}^l - \mu_{\theta}(x_t^l))^T \nabla_{\theta} \mu_{\theta}(x_t^l)] \right] \\ f_t &\triangleq \beta(1 - \sigma(\beta \log \frac{\pi_{\theta}(x_{t-1}^w | x_t^w, c)}{\pi_{\text{ref}}(x_{t-1}^w | x_t^w, c)} - \beta \log \frac{\pi_{\theta}(x_{t-1}^l | x_t^l, c)}{\pi_{\text{ref}}(x_{t-1}^l | x_t^l, c)})) \end{aligned} \quad (8)$$

In the above equation, the gradient is significantly affected by the relationship between inputs  $x_t^w$  and  $x_t^l$  from the previous step. This is because the input conditions  $(x_t^w, x_t^l)$  of generation probabilities for preferred sample  $x_{t-1}^w$  and dis-preferred sample  $x_{t-1}^l$  in Eq. (5) are different. Therefore, the choice of  $x_t^w$  and  $x_t^l$  disturbs the original optimization direction of DPO. In particular, if  $\nabla_{\theta} \mu_{\theta}(x_t^w) \approx \nabla_{\theta} \mu_{\theta}(x_t^l)$ , then the gradient term can be written as:

$$\nabla_{\theta} \mathcal{L}_{\text{D3PO}}(\theta) \approx -\mathbb{E} \left[ (f_t / \sigma_t^2) \cdot \nabla_{\theta}^T \mu_{\theta}(x_t^w) [(x_{t-1}^w - x_{t-1}^l) + (\mu_{\theta}(x_t^l) - \mu_{\theta}(x_t^w))] \right] \quad (9)$$

It shows that if  $x_t^w$  and  $x_t^l$  are located in the same linear subspace, then the optimization direction of the model shifts towards the direction  $\mu_{\theta}(x_t^l) - \mu_{\theta}(x_t^w)$ , which points to the dis-preferred samples. Thus, the fine-tuning effectiveness of DPO is significantly weakened.

### 3.3 TAILORED PREFERENCE OPTIMIZATION FRAMEWORK FOR DIFFUSION MODELS

To address the aforementioned problems, considering the characteristics of diffusion models, we propose a **Tailored Preference Optimization** (TailorPO) framework for fine-tuning diffusion models

in this section. Specifically, given a text prompt  $c$  and the time step  $t$ , we always start from the **same**  $x_t$  to generate the next time-step noisy samples, *i.e.*,  $x_{t-1}^0$  and  $x_{t-1}^1$ . Then, we estimate the step-wise reward of intermediate noisy samples  $x_{t-1}^0$  and  $x_{t-1}^1$  to directly rank their preference order. The sample with the higher reward value is represented by  $x_{t-1}^w$ , and the sample with the lower reward is denoted as  $x_{t-1}^l$ . Furthermore, if the reward function is differentiable, we apply the gradient guidance of the reward function (introduced in Section 3.4) to increase the reward of the preferred sample  $x_{t-1}^w$ , which enlarges the reward gap between  $x_{t-1}^w$  and  $x_{t-1}^l$  and enhances the fine-tuning effectiveness. At the next denoising step ( $t-1$ ), the preferred sample  $x_{t-1}^w$  is taken as  $x_{t-1}$  for further sampling and training. Our framework is illustrated in Figure 3, and the loss function is given as follows.

$$\mathcal{L}(\theta) = -\mathbb{E}_{(c, x_t, x_{t-1}^w, x_{t-1}^l)} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(x_{t-1}^w | x_t, c)}{\pi_{\text{ref}}(x_{t-1}^w | x_t, c)} - \beta \log \frac{\pi_\theta(x_{t-1}^l | x_t, c)}{\pi_{\text{ref}}(x_{t-1}^l | x_t, c)} \right) \right] \quad (10)$$

We will subsequently elucidate and substantiate the advantages of our proposed TailorPO framework for diffusion models from the following perspectives.

**Consistency between gradient direction and preferred samples.** First, TailorPO addresses the problem with the gradient direction of previous methods by always generating paired samples from the same  $x_t$ . This simple operation ensures that the generation probabilities used by the DPO loss function in Eq. (10) are all based on the same condition, aligning with the original formulation of DPO in Eq. (3). In this way, the gradient of our loss function is given as follows (please refer to Appendix A for the proof).

$$\nabla_\theta \mathcal{L}(\theta) = -\mathbb{E} \left[ (f_t / \sigma_t^2) \cdot \nabla_\theta^T \mu_\theta(x_t) (x_{t-1}^w - x_{t-1}^l) \right] \quad (11)$$

Notably, the gradient direction of our loss function clearly points towards the preferred samples. Therefore, the model is effectively encouraged to generate preferred samples.

**Immediate preference ranking at intermediate steps.** Instead of performing preference ranking on final images, we directly rank the preference order of noisy samples at intermediate steps. To this end, we propose to evaluate the preference quality of noisy samples  $x_t$ . As discussed in Section 3.2, the denoising process of a diffusion model can be formulated as an MDP, where the action value function for generating  $x_t$  simplifies to the expected reward of images over all trajectories starting from  $x_t$ . Therefore, we define the step-wise reward value of the noisy sample  $x_t$  as follows.

$$r_t(c, x_t) \triangleq \mathbb{E}[r(c, x_0) | c, x_t] \approx r(c, \hat{x}_0(x_t)) \quad (12)$$

However, computing the above expectation over all trajectories is intractable. Therefore, we employ an approximation to the expectation value. Previous studies (Chung et al., 2023; Guo et al., 2024) have proven that  $\mathbb{E}[x_0 | c, x_t] = \hat{x}_0(x_t)$ , which represents the predicted  $x_0$  at step  $t$  (defined in Eq. (2)). Furthermore, Chung et al. (2023) prove the following Proposition 1, which ensures that the expectation of image rewards  $\mathbb{E}[r(c, x_0) | c, x_t]$  can be approximated by the reward of the expected image  $r(c, \mathbb{E}[x_0 | c, x_t])$ . Therefore, we compute  $r_t(c, x_t) \approx r(c, \hat{x}_0(x_t))$  to estimate the step-wise reward of  $x_t$  for preference ranking.

**Proposition 1 (proven by Chung et al. (2023))** *Let a measurement  $g(x_0) = \mathcal{A}(x_0) + n$ , where  $\mathcal{A}(\cdot)$  is a measure operator defined on images  $x_0$  and  $n \sim \mathcal{N}(0, \sigma^2 I)$  is the measurement noise. The Jensen gap between  $\mathbb{E}[g(x_0) | c, x_t]$  and  $g(\mathbb{E}[x_0 | c, x_t])$ , *i.e.*,  $\mathcal{J} = \mathbb{E}[g(x_0) | c, x_t] - g(\mathbb{E}[x_0 | c, x_t])$  is bounded by  $\mathcal{J} \leq \frac{d}{\sqrt{2\pi\sigma^2}} e^{-1/2\sigma^2} \|\nabla_x \mathcal{A}(x)\|_{m_1}$ , where  $\nabla_x \mathcal{A}(x) \triangleq \max_x \|\nabla_x \mathcal{A}(x)\|$ ,  $m_1 \triangleq \int \|x_0 - \hat{x}_0\| p(x_0 | c, x_t) dx_0$ , and  $\hat{x}_0 = \mathbb{E}[x_0 | c, x_t]$ . The Jensen gap can approach 0 as  $\sigma$  increases.*

By obtaining the preference order of noisy samples immediately at intermediate steps, we can fine-tune the model using Eq. (10). Then, the preferred sample  $x_{t-1}^w$  is assigned as the input for the next step, enabling sampling and optimization in subsequent steps.

### 3.4 GRADIENT GUIDANCE OF REWARD MODEL FOR FINE-TUNING

In TailorPO, since noisy samples ( $x_{t-1}^0, x_{t-1}^1$ ) are generated from the same  $x_t$ , their similarity increases as  $t$  decreases. This increasing similarity potentially reduces the diversity of paired samples

Table 1: Gradient guidance successfully increased/decreased the reward of most samples.

$t$	20	16	12	8	4
ratio of $r_t(c, x_{t-1}^+) > r_t(c, x_{t-1})$	0.83	0.97	0.98	0.99	0.99
ratio of $r_t(c, x_{t-1}^-) < r_t(c, x_{t-1})$	0.87	0.98	1.00	0.98	1.00

**Algorithm 1:** The TailorPO-G framework for aligning diffusion models with human preference.

**Input:** Diffusion model  $\pi_\theta(\cdot)$ , reference model  $\pi_{\text{ref}}(\cdot)$ , reward model  $r(\cdot)$

```

1 Sample a text prompt  $c$ ;
2 Initialize  $x_T \sim \mathcal{N}(0, \mathbf{I})$ ;
3 for  $t = T, \dots, 1$  do
4   Sample  $x_{t-1}^0, x_{t-1}^1$  from  $\pi_\theta(\cdot|x_t, c)$ ;
5   Rank  $x_{t-1}^0$  and  $x_{t-1}^1$  based on their step-wise rewards to obtain  $x_{t-1}^w$  and  $x_{t-1}^l$ ;
6   Inject gradient guidance to compute  $x_{t-1}^+ = x_{t-1}^w - \eta_t \nabla_{x_{t-1}^w} (r_{\text{high}} - r_t(c, x_{t-1}^w))^2$ ;
7   if  $r_t(c, x_{t-1}^+) > r_t(c, x_{t-1}^w)$  then
8     |  $x_{t-1}^w \leftarrow x_{t-1}^+$ 
9   end
10  Optimize  $\pi_\theta(\cdot)$  using Eq. (10);
11   $x_{t-1} \leftarrow x_{t-1}^w$ ;
12 end

```

**Output:** The fine-tuned diffusion model  $\pi_\theta(\cdot)$ .

for training. On the other hand, Khaki et al. (2024) have shown that a large difference between paired samples is beneficial to the DPO effectiveness. Therefore, to enhance the DPO performance in this case, we propose enlarging the difference between two noisy samples from the reward perspective.

To this end, we consider how to adjust the reward of a noisy sample  $x_{t-1}$ . Similar to (Guo et al., 2024), we use  $r_{\text{high}}$  to represent an expected higher reward. Then, the gradient of the conditional score function is  $\nabla_{x_{t-1}} \log p(x_{t-1}|r_{\text{high}}) = \nabla \log p(x_{t-1}) + \nabla_{x_{t-1}} \log p(r_{\text{high}}|x_{t-1})$ , where the first term  $\nabla \log p(x_{t-1})$  is estimated by the diffusion model itself, and the second term is to be estimated by the guidance. Guo et al. (2024) further prove the following relationship for estimation.

$$\nabla_{x_{t-1}} \log p(r_{\text{high}}|x_{t-1}) \propto \nabla_{x_{t-1}} \log p(r_{\text{high}}|\hat{x}_0(x_{t-1})) \propto -\eta_t \nabla_{x_{t-1}} (r_{\text{high}} - r_t(c, x_{t-1}))^2 \quad (13)$$

Therefore, we can inject the gradient term  $\nabla_{x_{t-1}} (r_{\text{high}} - r_t(c, x_{t-1}))^2$  as the guidance to the generation of  $x_{t-1}$  to adjust its reward. Specifically, we update the noisy samples as follows.

$$\begin{aligned} x_{t-1}^+ &\leftarrow x_{t-1} - \eta_t \nabla_{x_{t-1}} (r_{\text{high}} - r_t(c, x_{t-1}))^2, \text{ to increase reward} \\ x_{t-1}^- &\leftarrow x_{t-1} + \eta_t \nabla_{x_{t-1}} (r_{\text{high}} - r_t(c, x_{t-1}))^2, \text{ to decrease reward} \end{aligned} \quad (14)$$

To demonstrate that the above gradient guidance is able to adjust the reward of noisy samples as expected, we compared the step-wise rewards of the original sample  $x_{t-1}$ , the increased sample  $x_{t-1}^+$ , and the decreased sample  $x_{t-1}^-$ . Specifically, we generated 100 noisy samples  $x_{t-1}$  from Stable Diffusion v1.5 (Rombach et al., 2022), and then computed the corresponding  $x_{t-1}^+$  and  $x_{t-1}^-$ . We set  $\eta_t = 0.2$  and  $r_{\text{high}} = r_t(c, x_{t-1}) + \delta$  following Guo et al. (2024), where the constant  $\delta = 0.5$  specified the expected increment of the reward value.

Then, we computed the ratio of increased samples (satisfying  $r_t(c, x_{t-1}^+) > r_t(c, x_{t-1})$ ) and the ratio of decreased samples (satisfying  $r_t(c, x_{t-1}^-) < r_t(c, x_{t-1})$ ). Table 1 shows that for almost all samples, the gradient guidance successfully increased or decreased their reward as expected, demonstrating its effectiveness in adapting the reward of samples.

Finally, we apply this method in our training process to enlarge the reward gap between a pair of noisy samples and develop the *TailorPO-G* framework. As shown in Figure 3 and Algorithm 1, we first modify the preferred sample  $x_{t-1}^w$  to increase its reward value, and then use the modified sample for fine-tuning and subsequent sampling.

## 4 EXPERIMENTS

**Experimental settings.** In our experiments, we evaluate the effectiveness of our method in fine-tuning Stable Diffusion v1.5 (Rombach et al., 2022). We compared our TailorPO method with the

Table 2: Reward values of images generated by diffusion models fine-tuned using different methods. The prompts are related to common animals.

	Aesthetic scorer	ImageReward	HPSv2	PickScore	Compressibility
Stable Diffusion v1.5	5.79	0.65	27.51	20.20	-105.51
DDPO (Black et al., 2024)	6.57	0.99	28.00	20.24	-37.37
D3PO (Yang et al., 2024a)	6.46	0.95	27.80	20.40	-29.31
SPO (Liang et al., 2024)	5.89	0.95	27.88	20.38	-
TailorPO	6.66	1.20	<b>28.37</b>	20.34	<b>-6.71</b>
TailorPO-G	<b>6.96</b>	<b>1.26</b>	28.03	<b>20.68</b>	-

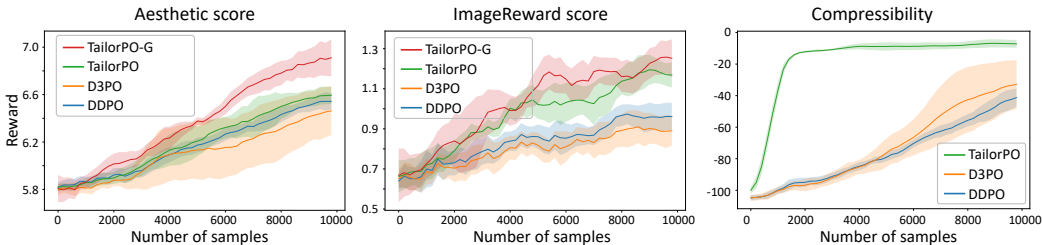


Figure 4: The change curve of reward values during the fine-tuning process. Experiments were conducted for three runs and we plot the average value and standard deviation of the reward.

RLHF method, DDPO (Black et al., 2024), and DPO-style methods, including D3PO (Yang et al., 2024a) and SPO (Liang et al., 2024). For all methods, we used the aesthetic scorer (Schuhmann et al., 2022), ImageReward (Xu et al., 2023), PickScore (Kirstain et al., 2023), HPSv2 (Wu et al., 2023), and JPEG compressibility measurement (Black et al., 2024) as reward models. Considering that some reward models are non-differentiable, we evaluate both the effectiveness of TailorPO and TailorPO-G, respectively.

Following the settings in D3PO (Yang et al., 2024a) and SPO (Liang et al., 2024), we applied the DDIM scheduler (Song et al., 2021) with  $\eta = 1.0$  and  $T = 20$  inference steps. The generated images were of resolution of  $512 \times 512$ . We employed LoRA (Hu et al., 2022) to fine-tune the UNet parameters on a total of 10,000 samples with a batch size of 2. The reference model was set as the pre-trained Stable Diffusion v1.5 itself. For SPO, we used the same hyper-parameters as in its original paper, and for other methods, we used the same hyper-parameters as in (Yang et al., 2024a), except that we set a smaller batch size. In particular, for all our frameworks, we generated images with  $T = 20$  and uniformly sampled  $T_{\text{fine-tune}} = 5$  steps for fine-tuning, *i.e.*, we only fine-tuned the model at steps  $t = 20, 16, 12, 8, 4$ . In addition, we set the coefficient  $\eta_t$  in gradient guidance using a cosine scheduler in the range of  $[0.1, 0.2]$ , which assigned a higher coefficient to smaller  $t$  (samples closer to output images). We have conducted ablation studies in Appendix C to show that our method is relatively stable with respect to the setting of  $T_{\text{fine-tune}}$  and  $\eta_t$ .

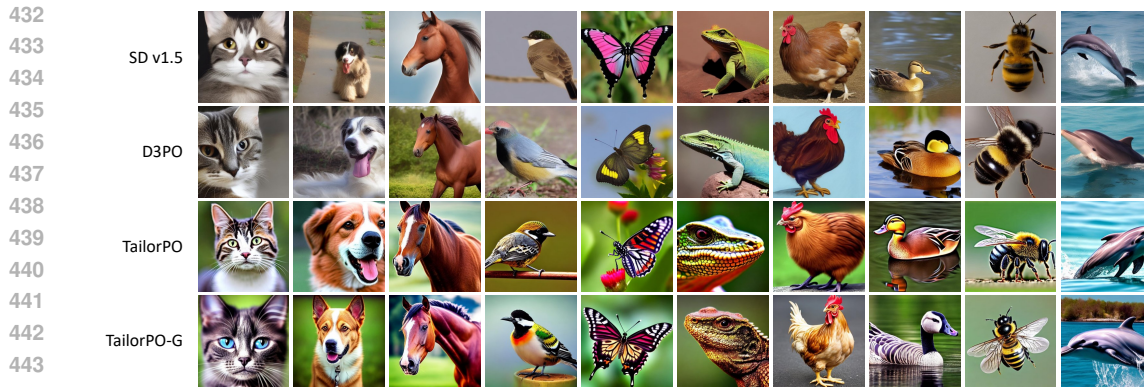
#### 4.1 EFFECTIVENESS OF ALIGNING DIFFUSION MODELS WITH PREFERENCE

In this section, we demonstrate that our frameworks outperform previous methods in aligning diffusion models with various preferences, from both quantitative and qualitative perspectives.

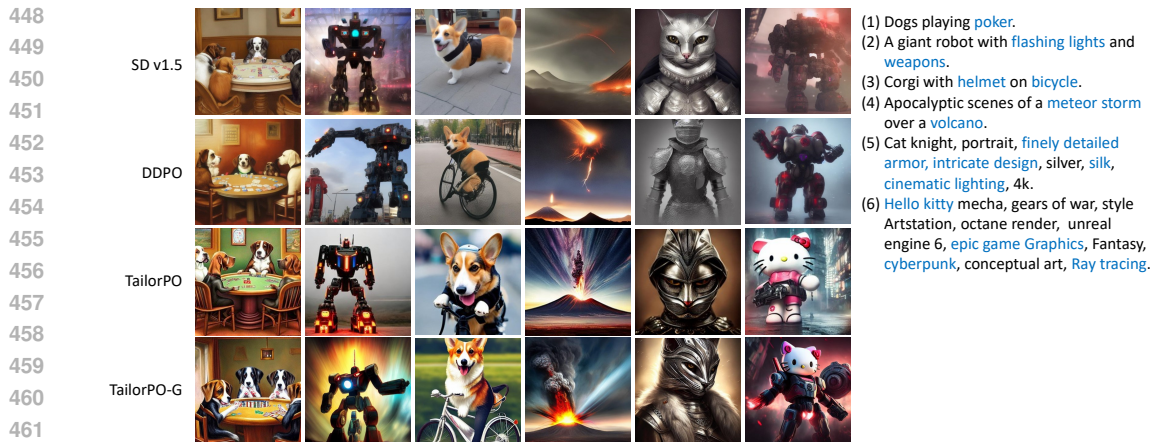
**Quantitative evaluation.** We fine-tuned Stable Diffusion v1.5 on various reward models using a set of prompts of common animals released by Black et al. (2024). For quantitative evaluation, we randomly sampled five images for each prompt and computed the average reward value of all images. Table 2 demonstrates that both TailorPO and TailorPO-G outperform other methods across all reward models. On the other hand, Figure 4 shows curves of reward values throughout the fine-tuning process. It can be observed that both of our frameworks rapidly increase the reward of generations in early iterations.

**Qualitative comparison.** For qualitative comparison, we first visualize the generated samples given simple prompts of animals in Figure 5. It is obvious that after fine-tuning using TailorPO and TailorPO-G, the model generated more colorful and visually appealing images with fine-grained details. In addition, we fine-tuned Stable Diffusion v1.5 on more complex prompts, using prompts in the Pick-a-Pic training dataset (Kirstain et al., 2023). Figure 6 shows that both TailorPO and TailorPO-G encourage the model to generate more aesthetically pleasing images, and these images





445 Figure 5: Visualization of images generated by diffusion models fine-tuned using different methods.  
446 For these animal-related prompts, diffusion models fine-tuned by TailorPO and TailorPO-G gener-  
447 ated more colorful and visually pleasing images.



463 Figure 6: Visualization of images generated by diffusion models fine-tuned on complex prompts in  
464 the Pick-a-Pic dataset. Prompts are given on the right with missing elements in SD v1.5 highlighted.  
465

466 were better aligned with the given prompts. For example, in the third row of Figure 6, the 5th and  
467 6th images contained more consistent and aligned subjects, scenes, and elements with the prompts.  
468

469 **User study.** Additionally, we conducted a user study by requesting five users to label their preference  
470 for generated images from the perspective of visual appeal and general preference. For each fine-  
471 tuned model, we generated five images for each animal-related prompt, Figure 7 reports the win-lose  
472 percentage results of our method versus other baseline methods, where our method exhibits a clear  
473 advantage in aligning with human preference. More experimental details can be seen in Appendix B.

474 4.2 GENERALIZATION TO DIFFERENT PROMPTS AND REWARD MODELS  
475

476 In this section, we investigate the generalization ability of the fine-tuned model using our method.  
477 Here, we consider two types of generalization mentioned in (Clark et al., 2024): prompt generaliza-  
478 tion and reward generalization.

479 **Prompt generalization** refers to the model’s ability to generate high-quality images for prompts  
480 beyond those used in fine-tuning. To evaluate this, we fine-tuned Stable Diffusion v1.5 on 45  
481 prompts of simple animal (Black et al., 2024) and evaluated its performance on 500 complex prompts  
482 (Kirstain et al., 2023). As shown in Table 3, the model fine-tuned on simple prompts exhibited  
483 higher reward values on complex prompts than the original SD v1.5, with our approach achieving  
484 the highest performance. Figure 8 presents examples of images generated from complex prompts,  
485 demonstrating that despite being fine-tuned on simple prompts, the model was also capable of gener-  
ating high-quality images given complex prompts. This highlights the effectiveness of our method

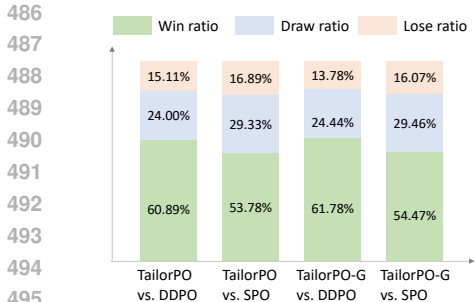


Figure 7: User-labeled win-lose ratio of TailorPO and TailorPO-G versus other baseline methods.



Figure 8: Diffusion model fine-tuned on simple prompts generalized well to complex prompts. Prompts from left to right are: (1) cinematic still of a stainless steel robot swimming in a pool. (2) A cat that is riding a horse without a leg. (3) Crazy frog, on one wheel, motorcycle, dead. (4) A panda riding a motorcycle. (5) Fantasy castle on a hilltop.

Table 3: Prompt generalization: the model fine-tuned on simple prompts also exhibited higher reward values for unseen complex prompts.

	Aesthetic scorer	ImageReward	HPSv2	PickScore	Compressibility
SD v1.5	5.69	-0.04	25.79	17.74	-98.95
DDPO	5.94	0.06	26.24	17.74	-49.94
D3PO	6.14	0.11	26.09	17.77	-38.92
SPO	5.79	0.15	26.28	17.16	-
TailorPO	6.26	0.11	<b>26.64</b>	17.85	<b>-7.32</b>
TailorPO-G	<b>6.45</b>	<b>0.25</b>	26.25	<b>17.93</b>	-

Table 4: Reward generalization: the model fine-tuned towards a reward model also exhibited higher reward values on other different but related reward models.

Train \ Evaluate	Aesthetic scorer	ImageReward	HPSv2	PickScore
SD v1.5	5.79	0.65	27.51	20.20
Aesthetic scorer	<b>6.96</b>	1.04	27.63	20.34
ImageReward	<u>6.01</u>	<b>1.26</b>	<u>28.01</u>	20.21
HPSv2	5.45	0.92	<b>28.03</b>	20.04
PickScore	5.94	0.83	27.71	<b>20.68</b>

in enhancing the model’s generalization to human-preferred images across various prompts, rather than overfitting to simple prompts.

**Reward generalization** refers to the phenomenon where fine-tuning the model towards a specific reward model can also enhance its performance on another different but related reward model. We selected one reward model from the aesthetic scorer, ImageReward, HPSv2, and PickScore for fine-tuning, and used the other three reward models for evaluation. Table 4 shows that after being fine-tuned towards the aesthetic scorer, ImageReward, and PickScore, the model usually exhibited higher performance on all these four reward models. In other words, our method boosted the overall ability of the model to generate high-quality images.

## 5 CONCLUSIONS

In this study, we rethink the existing DPO framework for aligning diffusion models and identify the potential flaws in these methods. We analyze these issues from both perspectives of preference order and gradient direction. To address these challenges, we consider the unique characteristics of diffusion models and introduce a novel tailored preference optimization framework for aligning diffusion models with human preference. Specifically, at each denoising step, our approach generates noisy samples from the same input and directly ranks their preference order for optimization. Furthermore, we propose integrating gradient guidance into the training framework to enhance the training effectiveness. Experimental results demonstrate that our approach significantly improved the reward scores of generated images, and exhibited good generalization over different prompts and different reward models.

## REFERENCES

- 540  
541  
542 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn  
543 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson  
544 Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernan-  
545 dez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson,  
546 Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and  
547 Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human  
548 feedback. *CoRR*, abs/2204.05862, 2022.
- 549 Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion  
550 models with reinforcement learning. In *ICLR*. OpenReview.net, 2024.
- 551 Chang Chen, Fei Deng, Kenji Kawaguchi, Caglar Gulcehre, and Sungjin Ahn. Simple hierarchical  
552 planning with diffusion. In *ICLR*. OpenReview.net, 2024.
- 553 Hyungjin Chung, Jeongsol Kim, Michael Thompson McCann, Marc Louis Klasky, and Jong Chul  
554 Ye. Diffusion posterior sampling for general noisy inverse problems. In *ICLR*. OpenReview.net,  
555 2023.
- 556 Kevin Clark, Paul Vicol, Kevin Swersky, and David J. Fleet. Directly fine-tuning diffusion models  
557 on differentiable rewards. In *ICLR*. OpenReview.net, 2024.
- 558 Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat gans on image synthesis. In  
559 *NeurIPS*, pp. 8780–8794, 2021.
- 560 Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,  
561 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for  
562 fine-tuning text-to-image diffusion models. In A. Oh, T. Naumann, A. Globerson, K. Saenko,  
563 M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36,  
564 pp. 79858–79885. Curran Associates, Inc., 2023.
- 565 Jiatao Gu, Alex Trevithick, Kai-En Lin, Joshua M. Susskind, Christian Theobalt, Lingjie Liu, and  
566 Ravi Ramamoorthi. Nerfdiff: Single-image view synthesis with nerf-guided distillation from  
567 3d-aware diffusion. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pp.  
568 11808–11826. PMLR, 2023.
- 569 Yingqing Guo, Hui Yuan, Yukang Yang, Minshuo Chen, and Mengdi Wang. Gradient guidance for  
570 diffusion models: An optimization perspective. *CoRR*, abs/2404.14743, 2024.
- 571 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *CoRR*, abs/2207.12598, 2022.
- 572 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*,  
573 2020.
- 574 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey A. Gritsenko,  
575 Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Im-  
576 agen video: High definition video generation with diffusion models. *CoRR*, abs/2210.02303,  
577 2022a.
- 578 Jonathan Ho, Chitwan Saharia, William Chan, David J. Fleet, Mohammad Norouzi, and Tim Sali-  
579 mans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*, 23:  
580 47:1–47:33, 2022b.
- 581 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
582 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR*. OpenRe-  
583 view.net, 2022.
- 584 Michael Janner, Yilun Du, Joshua B. Tenenbaum, and Sergey Levine. Planning with diffusion for  
585 flexible behavior synthesis. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*,  
586 pp. 9902–9915. PMLR, 2022.
- 587 Saeed Khaki, JinJin Li, Lan Ma, Liu Yang, and Prathap Ramachandra. RS-DPO: A hybrid rejection  
588 sampling and direct preference optimization method for alignment of large language models. In  
589 *NAACL-HLT (Findings)*, pp. 1665–1680. Association for Computational Linguistics, 2024.
- 590  
591  
592  
593

- 594 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-  
595 a-pic: An open dataset of user preferences for text-to-image generation. In *NeurIPS*, 2023.  
596
- 597 Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel,  
598 Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human  
599 feedback. *CoRR*, abs/2302.12192, 2023.
- 600 Zhanhao Liang, Yuhui Yuan, Shuyang Gu, Bohan Chen, Tiankai Hang, Ji Li, and Liang Zheng.  
601 Step-aware preference optimization: Aligning preference with denoising performance at each  
602 step. *CoRR*, abs/2406.04314, 2024.  
603
- 604 Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo P. Mandic, Wenwu Wang, and  
605 Mark D. Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. In *ICML*,  
606 volume 202 of *Proceedings of Machine Learning Research*, pp. 21450–21474. PMLR, 2023.
- 607 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob  
608 McGrew, Ilya Sutskever, and Mark Chen. GLIDE: towards photorealistic image generation and  
609 editing with text-guided diffusion models. In *ICML*, volume 162 of *Proceedings of Machine  
610 Learning Research*, pp. 16784–16804. PMLR, 2022.  
611
- 612 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023.
- 613 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,  
614 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser  
615 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan  
616 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback.  
617 In *NeurIPS*, 2022.  
618
- 619 Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d  
620 diffusion. In *ICLR*. OpenReview.net, 2023.
- 621 Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-  
622 image diffusion models with reward backpropagation. *CoRR*, abs/2310.03739, 2023.  
623
- 624 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and  
625 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.  
626 In *NeurIPS*, 2023.
- 627 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
628 resolution image synthesis with latent diffusion models. In *CVPR*, pp. 10674–10685. IEEE, 2022.  
629
- 630 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi  
631 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,  
632 Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jit-  
633 sev. LAION-5B: an open large-scale dataset for training next generation image-text models. In  
634 *NeurIPS*, 2022.
- 635 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry  
636 Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video:  
637 Text-to-video generation without text-video data. In *ICLR*. OpenReview.net, 2023.  
638
- 639 Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper-  
640 vised learning using nonequilibrium thermodynamics. In *ICML*, volume 37 of *JMLR Workshop  
641 and Conference Proceedings*, pp. 2256–2265. JMLR.org, 2015.
- 642 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*.  
643 OpenReview.net, 2021.  
644
- 645 Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,  
646 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using  
647 direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
and Pattern Recognition (CVPR)*, pp. 8228–8238, June 2024.

648 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.  
649 Human preference score v2: A solid benchmark for evaluating human preferences of text-to-  
650 image synthesis. *CoRR*, abs/2306.09341, 2023.  
651

652 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao  
653 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.  
654 In *NeurIPS*, 2023.

655 Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiabin Chen, Weihang Shen, Xiaolong Zhu, and Xiu Li.  
656 Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings of*  
657 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8941–8951,  
658 June 2024a.

659 Shentao Yang, Tianqi Chen, and Mingyuan Zhou. A dense reward view on aligning text-to-image  
660 diffusion with preference. In *ICML*. OpenReview.net, 2024b.  
661

662 Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten  
663 Kreis. LION: latent point diffusion models for 3d shape generation. In *NeurIPS*, 2022.  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A GRADIENT OF LOSS FUNCTIONS

**Gradient of the original DPO loss function.** Given the input  $(x, y^w, y^l) \sim \mathcal{D}$ , the loss of DPO is as follows.

$$\mathcal{L} = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(\beta \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)})] \quad (15)$$

Let  $h_\theta(x, y_w, y_l) \triangleq \beta \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)}$  and  $f(x, y_w, y_l) \triangleq \beta(1 - \sigma(h_\theta(x, y_w, y_l)))$ , then

$$\begin{aligned} \frac{\partial \mathcal{L}(x, y_w, y_l)}{\partial \theta} &= \frac{\partial -\log \sigma(h_\theta(x, y_w, y_l))}{\partial \theta} \\ &= -\frac{1}{\sigma(h_\theta(x, y_w, y_l))} \frac{\partial \sigma(h_\theta(x, y_w, y_l))}{\partial \theta} \\ &= -\frac{1}{\sigma(h_\theta(x, y_w, y_l))} \frac{\partial \sigma(h_\theta(x, y_w, y_l))}{\partial h_\theta(x, y_w, y_l)} \frac{\partial h_\theta(x, y_w, y_l)}{\partial \theta} \\ &= -\frac{1}{\sigma(h_\theta(x, y_w, y_l))} \sigma(h_\theta(x, y_w, y_l))(1 - \sigma(h_\theta(x, y_w, y_l))) \frac{\partial h_\theta(x, y_w, y_l)}{\partial \theta} \\ &= -f(x, y_w, y_l) \frac{\partial [\log \pi_\theta(y_w|x) - \log \pi_{\text{ref}}(y_w|x) - \log \pi_\theta(y_l|x) + \log \pi_{\text{ref}}(y_l|x)]}{\partial \theta} \\ &= -f(x, y_w, y_l) \left( \frac{\partial \log \pi_\theta(y_w|x)}{\partial \theta} - \frac{\partial \log \pi_\theta(y_l|x)}{\partial \theta} \right) \end{aligned} \quad (16)$$

**Gradient of the loss function of D3PO.** To study the generative distribution in the denoising process of diffusion models, let  $x \triangleq (x_t, c), y \triangleq x_{t-1}$ , then we have

$$\pi_\theta(y|x) = \pi_\theta(x_{t-1}|x_t, c) = \frac{1}{(2\pi\sigma_t^2)^{d/2}} \exp\left(-\frac{\|x_{t-1} - \mu_\theta(x_t)\|_2^2}{2\sigma_t^2}\right) \quad (17)$$

In this case, the gradient of the loglikelihood  $\log \pi_\theta(x_{t-1}|x_t, c)$  w.r.t.  $\theta$  is given as follows.

$$\begin{aligned} \frac{\partial \log \pi_\theta(x_{t-1}|x_t, c)}{\partial \theta} &= \left( \frac{\partial \mu_\theta(x_t)}{\partial \theta} \right)^T \frac{\partial \left( -\frac{\|x_{t-1} - \mu_\theta(x_t)\|_2^2}{2\sigma_t^2} - \log((2\pi\sigma_t^2)^{d/2}) \right)}{\partial \mu_\theta(x_t)} \\ &= \left( \frac{\partial \mu_\theta(x_t)}{\partial \theta} \right)^T \frac{(x_{t-1} - \mu_\theta(x_t))}{\sigma_t^2} \end{aligned} \quad (18)$$

Then, we consider the gradient of the D3PO loss w.r.t. the model output  $\mu_\theta$ .

$$\begin{aligned} \frac{\partial \mathcal{L}(x_t^w, x_{t-1}^w, x_t^l, x_{t-1}^l)}{\partial \theta} &= -f_t \left( \frac{\partial \log \pi_\theta(x_{t-1}^w|x_t^w, t, c)}{\partial \theta} - \frac{\partial \log \pi_\theta(x_{t-1}^l|x_t^l, t, c)}{\partial \theta} \right) \\ &= -\frac{f_t}{\sigma_t^2} \left[ \left( \frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \right)^T (x_{t-1}^w - \mu_\theta(x_t^w)) - \left( \frac{\partial \mu_\theta(x_t^l)}{\partial \theta} \right)^T (x_{t-1}^l - \mu_\theta(x_t^l)) \right] \end{aligned} \quad (19)$$

Suppose  $\Delta\theta = -\eta \frac{\partial \mathcal{L}(x_t^w, x_{t-1}^w, x_t^l, x_{t-1}^l)}{\partial \theta}$ . After the update of  $\theta' \leftarrow \theta + \Delta\theta$ ,  $\Delta\mu_\theta(x_t^w) \approx \eta \frac{f_t}{\sigma_t^2} \left[ \left( \frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \right) \left( \frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \right)^T (x_{t-1}^w - \mu_\theta(x_t^w)) \right] - \eta \frac{f_t}{\sigma_t^2} \left[ \left( \frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \right) \left( \frac{\partial \mu_\theta(x_t^l)}{\partial \theta} \right)^T (x_{t-1}^l - \mu_\theta(x_t^l)) \right]$ . If  $x_t^w$  and  $x_t^l$  are located in the same linear subspace of the model, i.e.,  $\frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \approx \frac{\partial \mu_\theta(x_t^l)}{\partial \theta}$ , then the gradient can be written as follows.

$$\begin{aligned} \frac{\partial \mathcal{L}(x_t^w, x_{t-1}^w, x_t^l, x_{t-1}^l)}{\partial \theta} &= -\frac{f_t}{\sigma_t^2} \left[ \left( \frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \right)^T (x_{t-1}^w - \mu_\theta(x_t^w)) - \left( \frac{\partial \mu_\theta(x_t^l)}{\partial \theta} \right)^T (x_{t-1}^l - \mu_\theta(x_t^l)) \right] \\ &\approx -\frac{f_t}{\sigma_t^2} \left[ \left( \frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \right)^T (x_{t-1}^w - \mu_\theta(x_t^w)) - \left( \frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \right)^T (x_{t-1}^l - \mu_\theta(x_t^l)) \right] \\ &\approx -\frac{f_t}{\sigma_t^2} \left( \frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \right)^T [(x_{t-1}^w - x_{t-1}^l) + (\mu_\theta(x_t^l) - \mu_\theta(x_t^w))] \end{aligned} \quad (20)$$

Table 5: Effect of the number of steps used in TailorPO. For each setting of  $T_{\text{fine-tune}}$ , we uniformly sampled  $T_{\text{fine-tune}}$  steps for fine-tuning.

$T_{\text{fine-tune}}$	Aesthetic Scorer	HPSv2	compressibility
10	6.61	28.14	-20.62
5	<b>6.74</b>	<b>28.43</b>	<b>-4.76</b>
3	6.40	28.15	-9.97

Table 6: Effect of strength  $\eta_t$  of gradient guidance in TailorPO-G. [0.1,0.2] represents we set  $\eta_t$  ranging from 0.1 to 0.2 for different  $t$ .

$\eta_t$	Aesthetic Scorer	ImageReward	HPSv2
0.1	5.82	1.22	28.10
0.2	6.97	<b>1.35</b>	28.18
0.5	7.07	0.71	27.48
[0.1, 0.2]	<b>7.11</b>	1.25	<b>28.43</b>

Suppose  $\Delta\theta = -\eta \frac{\partial \mathcal{L}(x_t^w, x_{t-1}^w, x_t^l, x_{t-1}^l)}{\partial \theta}$ . After the update of  $\theta' \leftarrow \theta + \Delta\theta$ ,  $\Delta\mu_\theta(x_t^w) \approx \eta \frac{f_t}{\sigma_t^2} \left( \frac{\partial \mu_\theta(x_t^w)}{\partial \theta} \right) \left( \frac{\partial \mu_\theta(x_t^l)}{\partial \theta} \right)^T [(x_{t-1}^w - x_{t-1}^l) + (\mu_\theta(x_t^l) - \mu_\theta(x_t^w))]$ .

**Gradient of our loss function.** Then, we consider the gradient of our loss function *w.r.t.* the model output  $\mu_\theta$ .

$$\begin{aligned} \frac{\partial \mathcal{L}(x_t, x_{t-1}^w, x_{t-1}^l)}{\partial \theta} &= -f_t \left( \frac{\partial \mu_\theta(x_t)}{\partial \theta} \right)^T \left( \frac{\partial \log \pi_\theta(x_{t-1}^w | x_t, t, c)}{\partial \mu_\theta(x_t)} - \frac{\partial \log \pi_\theta(x_{t-1}^l | x_t, t, c)}{\partial \mu_\theta(x_t)} \right) \\ &= -f_t \left( \frac{\partial \mu_\theta(x_t)}{\partial \theta} \right)^T \left( \frac{x_{t-1}^w - \mu_\theta(x_t)}{\sigma_t^2} - \frac{x_{t-1}^l - \mu_\theta(x_t)}{\sigma_t^2} \right) \\ &= -\frac{f_t}{\sigma_t^2} \left( \frac{\partial \mu_\theta(x_t)}{\partial \theta} \right)^T (x_{t-1}^w - x_{t-1}^l) \end{aligned} \quad (21)$$

Suppose  $\Delta\theta = -\eta \frac{\partial \mathcal{L}(x_t, x_{t-1}^w, x_{t-1}^l)}{\partial \theta}$ . After the update of  $\theta' \leftarrow \theta + \Delta\theta$ ,  $\Delta\mu_\theta(x_t) \approx \left( \frac{\partial \mu_\theta(x_t)}{\partial \theta} \right) \Delta\theta = \eta \frac{f_t}{\sigma_t^2} \left( \frac{\partial \mu_\theta(x_t)}{\partial \theta} \right) \left( \frac{\partial \mu_\theta(x_t)}{\partial \theta} \right)^T (x_{t-1}^w - x_{t-1}^l)$ .

## B EXPERIMENTAL SETTINGS FOR THE USER STUDY

To verify that our framework generates more human-preferred images, we conducted a user study by requesting five human users to label their preference for generated images from the perspective of visual appeal and general preference. Given each prompt in the set of 45 animal prompts, we sampled five images from the fine-tuned model and obtained a total of 225 images per model. For comparison, for each pair of fine-tuned model, we organized their generated images into 225 pairs. Users were then asked to compare each pair of images and label their preferences. If the images in a pair looked very similar or were both unappealing, the user labeled “draw” for them. Then, we computed the ratio of pairs where TailorPO and TailorPO-G received “win”. “draw”, and “lose” labels, respectively. Figure 7 reports the win-lose percentage results of our method versus other baseline methods, our method exhibits a clear advantage in aligning with human preference.

## C ABLATION STUDIES

In this section, we performed ablation studies to verify the effect of hyper-parameters on performance, including the number of steps used for optimization and the strength of gradient guidance.

**Effect of steps used for training.** We first investigate the effect of the number of steps  $T_{\text{fine-tune}}$  used for fine-tuning in TailorPO. In Section 4, We generated images with  $T = 20$  sampling timesteps and uniformly sampled only  $T_{\text{fine-tune}} = 5$  steps for training to boost the training efficiency. Here, we compared the results of setting  $T_{\text{fine-tune}} = 3, 5, 10$  in Table 5, and it shows that while the fine-tuning performance is relatively stable to the setting of  $T_{\text{fine-tune}}$ , fine-tuning on five steps achieved a better trade-off between performance and efficiency.

**Effect of the strength of gradient guidance.** We also verify the effect of gradient guidance in TailorPO-G by applying gradient guidance with different strengths at intermediate steps. Specifically, we used different settings of  $\eta_t$  in Eq. (14) for fine-tuning. The result in Table 6 shows that the varying strength  $\eta_t$  for different steps  $t$  better enhance the fine-tuning performance.