

# PRDP: Proximal Reward Difference Prediction for Large-Scale Reward Finetuning of Diffusion Models

## Supplementary Material

### A. Proofs

#### A.1. Lower Bound of RLHF Objective

In Lemma A.1, we prove that the objective in Equation (6) is a lower bound of the RLHF objective in Equation (5).

**Lemma A.1.** *Given two diffusion models  $\pi_\theta, \pi_{\text{ref}}$ , a prompt distribution  $p(\mathbf{c})$ , a reward function  $r(\mathbf{x}_0, \mathbf{c})$ , and a constant  $\beta > 0$ , we have:*

$$\mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \left[ \mathbb{E}_{\mathbf{x}_0 \sim \pi_\theta(\mathbf{x}_0 | \mathbf{c})} [r(\mathbf{x}_0, \mathbf{c})] - \beta \text{KL}[\pi_\theta(\mathbf{x}_0 | \mathbf{c}) || \pi_{\text{ref}}(\mathbf{x}_0 | \mathbf{c})] \right] \quad (22)$$

$$\geq \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \left[ \mathbb{E}_{\mathbf{x}_0 \sim \pi_\theta(\mathbf{x}_0 | \mathbf{c})} [r(\mathbf{x}_0, \mathbf{c})] - \beta \text{KL}[\pi_\theta(\bar{\mathbf{x}} | \mathbf{c}) || \pi_{\text{ref}}(\bar{\mathbf{x}} | \mathbf{c})] \right], \quad (23)$$

where  $\bar{\mathbf{x}} := \mathbf{x}_{0:T}$  is the full denoising trajectory, and  $\pi_\theta, \pi_{\text{ref}}$  are defined as:

$$\pi(\mathbf{x}_0 | \mathbf{c}) = \int \pi(\mathbf{x}_{0:T} | \mathbf{c}) \, d\mathbf{x}_{1:T} = \int p(\mathbf{x}_T) \prod_{t=1}^T \pi(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{c}) \, d\mathbf{x}_{1:T}. \quad (24)$$

*Proof.* It suffices to show that for any  $\mathbf{c}$ ,

$$\text{KL}[\pi_\theta(\bar{\mathbf{x}} | \mathbf{c}) || \pi_{\text{ref}}(\bar{\mathbf{x}} | \mathbf{c})] \geq \text{KL}[\pi_\theta(\mathbf{x}_0 | \mathbf{c}) || \pi_{\text{ref}}(\mathbf{x}_0 | \mathbf{c})]. \quad (25)$$

This can be proved similarly as the data processing inequality. We provide the proof below.

$$\text{KL}[\pi_\theta(\bar{\mathbf{x}} | \mathbf{c}) || \pi_{\text{ref}}(\bar{\mathbf{x}} | \mathbf{c})] = \mathbb{E}_{\pi_\theta(\mathbf{x}_{0:T} | \mathbf{c})} \left[ \log \frac{\pi_\theta(\mathbf{x}_{0:T} | \mathbf{c})}{\pi_{\text{ref}}(\mathbf{x}_{0:T} | \mathbf{c})} \right] \quad (26)$$

$$= \mathbb{E}_{\pi_\theta(\mathbf{x}_{0:T} | \mathbf{c})} \left[ \log \frac{\pi_\theta(\mathbf{x}_0 | \mathbf{c})}{\pi_{\text{ref}}(\mathbf{x}_0 | \mathbf{c})} + \log \frac{\pi_\theta(\mathbf{x}_{1:T} | \mathbf{x}_0, \mathbf{c})}{\pi_{\text{ref}}(\mathbf{x}_{1:T} | \mathbf{x}_0, \mathbf{c})} \right] \quad (27)$$

$$= \mathbb{E}_{\pi_\theta(\mathbf{x}_0 | \mathbf{c})} \left[ \log \frac{\pi_\theta(\mathbf{x}_0 | \mathbf{c})}{\pi_{\text{ref}}(\mathbf{x}_0 | \mathbf{c})} \right] + \mathbb{E}_{\pi_\theta(\mathbf{x}_0 | \mathbf{c})} \left[ \mathbb{E}_{\pi_\theta(\mathbf{x}_{1:T} | \mathbf{x}_0, \mathbf{c})} \left[ \log \frac{\pi_\theta(\mathbf{x}_{1:T} | \mathbf{x}_0, \mathbf{c})}{\pi_{\text{ref}}(\mathbf{x}_{1:T} | \mathbf{x}_0, \mathbf{c})} \right] \right] \quad (28)$$

$$= \text{KL}[\pi_\theta(\mathbf{x}_0 | \mathbf{c}) || \pi_{\text{ref}}(\mathbf{x}_0 | \mathbf{c})] + \mathbb{E}_{\pi_\theta(\mathbf{x}_0 | \mathbf{c})} [\text{KL}[\pi_\theta(\mathbf{x}_{1:T} | \mathbf{x}_0, \mathbf{c}) || \pi_{\text{ref}}(\mathbf{x}_{1:T} | \mathbf{x}_0, \mathbf{c})]] \quad (29)$$

$$\geq \text{KL}[\pi_\theta(\mathbf{x}_0 | \mathbf{c}) || \pi_{\text{ref}}(\mathbf{x}_0 | \mathbf{c})]. \quad (30)$$

□

## A.2. Maximizer of the Lower Bound of RLHF Objective

In Lemma A.2, we prove that Equation (7) maximizes the objective in Equation (6), a lower bound of the RLHF objective.

**Lemma A.2.** *Define*

$$\pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c}) = \frac{1}{Z(\mathbf{c})} \pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c}) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, \mathbf{c})\right), \quad (31)$$

where

$$Z(\mathbf{c}) = \int \pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c}) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, \mathbf{c})\right) d\bar{\mathbf{x}} \quad (32)$$

is the partition function. Then  $\pi_{\theta^*}$  is the optimal solution to the following maximization problem:

$$\max_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \left[ \mathbb{E}_{\mathbf{x}_0 \sim \pi_{\theta}(\mathbf{x}_0|\mathbf{c})} [r(\mathbf{x}_0, \mathbf{c})] - \beta \text{KL}[\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c}) || \pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c})] \right]. \quad (33)$$

*Proof.* We provide the proof below, which is inspired by DPO [35].

$$\max_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \left[ \mathbb{E}_{\mathbf{x}_0 \sim \pi_{\theta}(\mathbf{x}_0|\mathbf{c})} [r(\mathbf{x}_0, \mathbf{c})] - \beta \text{KL}[\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c}) || \pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c})] \right] \quad (34)$$

$$= \max_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \left[ \mathbb{E}_{\bar{\mathbf{x}} \sim \pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})} [r(\mathbf{x}_0, \mathbf{c})] - \beta \text{KL}[\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c}) || \pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c})] \right] \quad (35)$$

$$= \max_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \mathbb{E}_{\bar{\mathbf{x}} \sim \pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})} \left[ r(\mathbf{x}_0, \mathbf{c}) - \beta \log \frac{\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})}{\pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c})} \right] \quad (36)$$

$$= \min_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \mathbb{E}_{\bar{\mathbf{x}} \sim \pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})} \left[ \log \frac{\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})}{\pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c})} - \frac{1}{\beta} r(\mathbf{x}_0, \mathbf{c}) \right] \quad (37)$$

$$= \min_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \mathbb{E}_{\bar{\mathbf{x}} \sim \pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})} \left[ \log \frac{\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})}{\pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c}) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, \mathbf{c})\right)} \right] \quad (38)$$

$$= \min_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \mathbb{E}_{\bar{\mathbf{x}} \sim \pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})} \left[ \log \frac{\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})}{\pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c}) Z(\mathbf{c})} \right] \quad (39)$$

$$= \min_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} \left[ \mathbb{E}_{\bar{\mathbf{x}} \sim \pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})} \left[ \log \frac{\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c})}{\pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c})} \right] - \log Z(\mathbf{c}) \right] \quad (40)$$

$$= \min_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} [\text{KL}[\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c}) || \pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c})] - \log Z(\mathbf{c})] \quad (41)$$

$$= \min_{\pi_{\theta}} \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c})} [\text{KL}[\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c}) || \pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c})]]. \quad (42)$$

Since  $\text{KL}[\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c}) || \pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c})] \geq 0$ , and  $\text{KL}[\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c}) || \pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c})] = 0$  if and only if  $\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c}) = \pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c})$ , we conclude that the optimal solution to Equation (33) is  $\pi_{\theta}(\bar{\mathbf{x}}|\mathbf{c}) = \pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c})$  for all  $\mathbf{c}$ .  $\square$

### A.3. Necessary and Sufficient Conditions for the Optimal Solution

In Lemma A.3, we provide theoretical justification for our proposed RDP objective in Equation (14).

**Lemma A.3.**

$$\pi_\theta(\bar{\mathbf{x}}|\mathbf{c}) = \pi_{\theta^*}(\bar{\mathbf{x}}|\mathbf{c}), \quad \forall \bar{\mathbf{x}}, \mathbf{c} \quad (43)$$

$$\iff \log \frac{\pi_\theta(\bar{\mathbf{x}}^a|\mathbf{c})}{\pi_{\text{ref}}(\bar{\mathbf{x}}^a|\mathbf{c})} - \log \frac{\pi_\theta(\bar{\mathbf{x}}^b|\mathbf{c})}{\pi_{\text{ref}}(\bar{\mathbf{x}}^b|\mathbf{c})} = \frac{r(\mathbf{x}_0^a, \mathbf{c}) - r(\mathbf{x}_0^b, \mathbf{c})}{\beta}, \quad \forall \bar{\mathbf{x}}^a, \bar{\mathbf{x}}^b, \mathbf{c}. \quad (44)$$

*Proof.* We have shown “ $\implies$ ” in the main text. We provide the proof for “ $\impliedby$ ” below.

Equation (44) implies that

$$\log \frac{\pi_\theta(\bar{\mathbf{x}}|\mathbf{c})}{\pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c})} - \frac{1}{\beta} r(\mathbf{x}_0, \mathbf{c}) \quad (45)$$

is a constant w.r.t.  $\bar{\mathbf{x}}$ . Therefore, we can write Equation (45) as a function of  $\mathbf{c}$  alone:

$$\log \frac{\pi_\theta(\bar{\mathbf{x}}|\mathbf{c})}{\pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c})} - \frac{1}{\beta} r(\mathbf{x}_0, \mathbf{c}) = f(\mathbf{c}). \quad (46)$$

Hence,

$$\pi_\theta(\bar{\mathbf{x}}|\mathbf{c}) = \pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c}) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, \mathbf{c})\right) \exp(f(\mathbf{c})). \quad (47)$$

It suffices to show that

$$\exp(f(\mathbf{c})) = \frac{1}{Z(\mathbf{c})}, \quad \forall \mathbf{c}. \quad (48)$$

This follows from the fact that the probability density function  $\pi_\theta(\bar{\mathbf{x}}|\mathbf{c})$  must satisfy:

$$1 = \int \pi_\theta(\bar{\mathbf{x}}|\mathbf{c}) d\bar{\mathbf{x}} \quad (49)$$

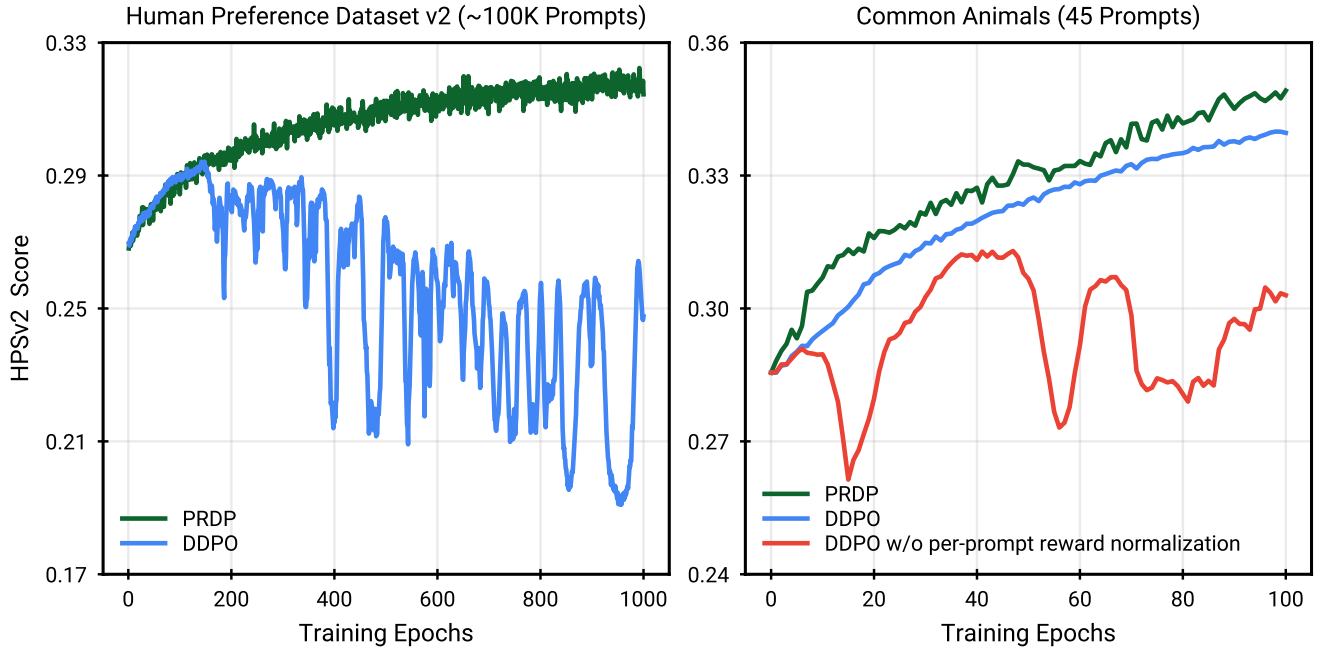
$$= \int \pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c}) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, \mathbf{c})\right) \exp(f(\mathbf{c})) d\bar{\mathbf{x}} \quad (50)$$

$$= \exp(f(\mathbf{c})) \int \pi_{\text{ref}}(\bar{\mathbf{x}}|\mathbf{c}) \exp\left(\frac{1}{\beta} r(\mathbf{x}_0, \mathbf{c})\right) d\bar{\mathbf{x}} \quad (51)$$

$$= \exp(f(\mathbf{c})) Z(\mathbf{c}). \quad (52)$$

□

## B. Instability of DDPO in Large-Scale Reward Finetuning

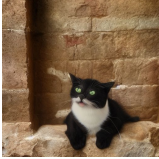
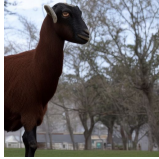

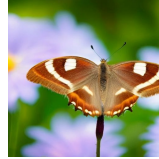
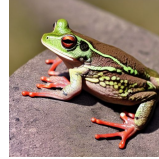






















**Figure 9. Analysis of the instability of DDPO in large-scale training.** We plot the training curves of PRDP and DDPO on the large-scale Human Preference Dataset v2 (Left) and the small-scale Common Animals (Right). PRDP outperforms DDPO in the small-scale setting, and maintains stability in the large-scale setting where DDPO fails. Our ablation study suggests that the per-prompt reward normalization in DDPO is key to its stability, and the inability to perform such normalization in the large-scale setting likely causes its failure.

Figure 9 shows the training curve of PRDP and DDPO [4], where the reward model is HPSv2 [53]. From Figure 9 (Left), we observe that when trained on the large-scale Human Preference Dataset v2 (HPD v2) [53], DDPO fails to stably optimize the reward. We conjecture that this is because the per-prompt reward normalization is rarely enabled in the large-scale setting, since each prompt can only be seen a few times. Specifically, in each epoch, DDPO randomly samples 512 prompts, so on average, each prompt can be seen  $512 \times 1000 / 100K \approx 5$  times. This is insufficient to obtain a good estimate of the per-prompt expected reward. In this case, DDPO will compute a prompt-agnostic expected reward, by averaging the rewards across all 512 prompts. To verify that such prompt-agnostic reward normalization causes training instability, we conduct an ablation study of DDPO in our small-scale setting with 45 training prompts. As shown in Figure 9 (Right), DDPO without per-prompt reward normalization is unstable even in the small-scale setting, suggesting that the inability to perform per-prompt reward normalization can be a limiting factor in scaling DDPO to large prompt datasets. In contrast to DDPO, PRDP can steadily improve the reward score and maintain stability in both small-scale and large-scale settings.



### C. Effect of KL Regularization

	cat	goat	tiger	butterfly	frog	
Stable Diffusion						HPSv2: 0.2836 PickScore: 0.2155 Aesthetic: 5.49
DDPO						HPSv2: 0.2629 PickScore: 0.1989 Aesthetic: 9.52
PRDP (beta = 0.1)						HPSv2: 0.2647 PickScore: 0.2038 Aesthetic: <b>9.70</b>
PRDP (beta = 1)						HPSv2: 0.2794 PickScore: 0.2199 Aesthetic: 8.62
PRDP (beta = 10)						HPSv2: <b>0.2841</b> PickScore: <b>0.2212</b> Aesthetic: 7.45

**Figure 10. Effect of KL regularization on optimizing aesthetic score.** DDPO and PRDP are finetuned from Stable Diffusion v1.4 on 45 prompts of common animal names. Evaluation is performed on the same set of prompts. In addition to aesthetic score, we report HPSv2 and PickScore which reflect text-image alignment but are not used during training. Samples within each column are generated from the prompt shown on top, using the same random seed. PRDP with a large KL weight  $\beta$  can alleviate the reward over-optimization problem encountered by DDPO, significantly improving the aesthetic quality over Stable Diffusion while maintaining text-image alignment.

In contrast to DDPO [4] which only cares about maximizing the reward, PRDP is formulated with a KL regularization, allowing us to alleviate the problem of reward over-optimization by increasing the KL weight  $\beta$ . We demonstrate the effect of KL regularization in Figure 10. Here, the reward used for training is the aesthetic score given by the LAION aesthetic predictor. It only takes images as input, and therefore ignores the text-image alignment. We finetune DDPO and PRDP from Stable Diffusion v1.4 [37] for 250 epochs on 45 training prompts of common animal names as used in DDPO, with 512 reward queries in each epoch. For evaluation, we additionally use HPSv2 [53] and PickScore [22] that reflect text-image alignment. The reported reward scores are averaged over 64 random samples per training prompt, using the same random seed for Stable Diffusion v1.4, DDPO, and PRDP.

We observe that DDPO, without KL regularization, is prone to reward over-optimization. It ignores the text prompt and generates similar images for all prompts. PRDP with a small KL weight (*e.g.*,  $\beta = 0.1$ ) has the same problem, but achieves higher reward scores than DDPO, showing a better reward maximization capability. As the KL weight increases, PRDP is able to better preserve the text-image alignment, indicated by the increase in HPSv2 and PickScore. With  $\beta = 10$ , PRDP significantly improves the aesthetic score over Stable Diffusion v1.4 without sacrificing text-image alignment.

## D. Large-Scale Multi-Reward Finetuning

**Table 3. Reward score comparison on unseen prompts.** We use a weighted combination of rewards: PickScore = 10, HPSv2 = 2, Aesthetic = 0.05. PRDP is finetuned from Stable Diffusion v1.4 on the training set prompts of Pick-a-Pic v1 dataset.

	Pick-a-Pic v1 Test Set	HPD v2 Animation	HPD v2 Concept Art	HPD v2 Painting	HPD v2 Photo
SD v1.4	2.888	2.927	2.877	2.883	2.984
PRDP	<b>3.208</b>	<b>3.296</b>	<b>3.264</b>	<b>3.274</b>	<b>3.214</b>

In this section, we provide additional results for our large-scale multi-reward finetuning experiment. Following DRaFT [6], we use a weighted combination of rewards: PickScore = 10, HPSv2 = 2, Aesthetic = 0.05. We finetune Stable Diffusion v1.4 [37] on the training set prompts of Pick-a-Pic v1 dataset [22]. We evaluate our finetuned model on a variety of unseen prompts, including 500 prompts from the Pick-a-Pic v1 test set, and 800 prompts from each of the four benchmark categories of the Human Preference Dataset v2 (HPD v2) [53], namely animation, concept art, painting, and photo. Table 3 reports the reward scores before and after finetuning. The reward scores are averaged over 64 random samples per prompt, using the same random seed for Stable Diffusion v1.4 and PRDP. We further show generation samples for each test prompt set in Figures 11 to 15. As can be seen, PRDP significantly improves generation quality across all five prompt sets.

## E. Hyperparameters

**Table 4. PRDP training hyperparameters.**

Name	Symbol	Small-Scale Finetuning	Large-Scale Finetuning	Large-Scale Multi-Reward Finetuning
Training epochs	$E$	100	1000	1000
Gradient updates per epoch	$K$	10	1	1
Prompts per epoch	$N$	32	64	64
Images per prompt	$B$	16	8	8
KL weight	$\beta$	$3 \times 10^{-5}$	$3 \times 10^{-6}$	$3 \times 10^{-5}$
DDPM steps	$T$	50	50	50
Stepwise clipping range	$\epsilon$	$1 \times 10^{-6}$	$1 \times 10^{-4}$	$1 \times 10^{-4}$
Classifier-free guidance scale	—	5.0	5.0	5.0
Optimizer	—	AdamW	AdamW	AdamW
Gradient clipping	—	1.0	1.0	1.0
Learning rate	—	$1 \times 10^{-5}$	$7 \times 10^{-6}$	$1 \times 10^{-5}$
Weight decay	—	$1 \times 10^{-4}$	$1 \times 10^{-4}$	$1 \times 10^{-4}$

## F. Effect of Clipping

Table 5. Effect of clipping on training stability.














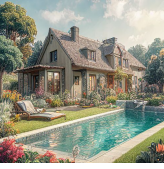







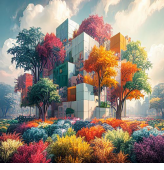




















	w/o Clipping	w/ Clipping
DDPO	Small scale: <b>Unstable</b> Large scale: <b>Unstable</b>	Small scale: <b>Stable</b> Large scale: <b>Unstable</b>
PRDP	Small scale: <b>Unstable</b> Large scale: <b>Unstable</b>	Small scale: <b>Stable</b> Large scale: <b>Stable</b>

Table 5 summarizes the effect of clipping on the training stability of both DDPO [4] and PRDP. For DDPO, we use PPO-based clipping [42], while for PRDP, we use the proximal updates described in Section 3.3. We observe that clipping is key to stability of small-scale training, whereas using the PRDP objective and clipping are both indispensable for achieving stability in large-scale training.

## G. Jax Implementation of PRDP Loss










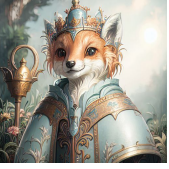
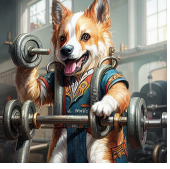


















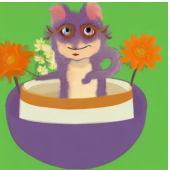


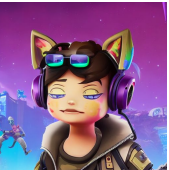









```
1 import jax
2 import jax.numpy as jnp
3
4
5 def prdp_loss(
6     log_probs: jax.Array,      # (B, T)
7     log_probs_old: jax.Array,  # (B, T)
8     log_probs_ref: jax.Array,  # (B, T)
9     rewards: jax.Array,       # (B,)
10    clip_range: float,
11    kl_weight: float,
12 ) -> jax.Array:
13     """Computes PRDP loss for a batch of denoising trajectories with the same text prompt.
14
15     Args:
16         log_probs: Log probs of the denoising trajectories under pi_theta.
17         log_probs_old: Log probs of the denoising trajectories under pi_theta_old.
18         log_probs_ref: Log probs of the denoising trajectories under pi_ref.
19         rewards: Rewards of the generated clean images.
20         clip_range: Stepwise clipping range (epsilon).
21         kl_weight: KL weight (beta).
22
23     Returns:
24         loss: The PRDP loss.
25     """
26     log_ratios = log_probs - log_probs_ref
27     log_ratios_old = log_probs_old - log_probs_ref
28     clipped_log_ratios = jnp.clip(
29         log_ratios, log_ratios_old - clip_range, log_ratios_old + clip_range
30     )
31
32     log_ratios = jnp.mean(log_ratios, axis=-1)
33     clipped_log_ratios = jnp.mean(clipped_log_ratios, axis=-1)
34
35     log_ratio_diffs = log_ratios[:, None] - log_ratios
36     clipped_log_ratio_diffs = clipped_log_ratios[:, None] - clipped_log_ratios
37     reward_diffs = rewards[:, None] - rewards
38
39     mse_loss = (log_ratio_diffs - reward_diffs / kl_weight) ** 2
40     clipped_mse_loss = (clipped_log_ratio_diffs - reward_diffs / kl_weight) ** 2
41     loss = jnp.maximum(mse_loss, clipped_mse_loss)
42     loss = jnp.mean(loss, where=reward_diffs > 0)
43
44     return loss
```



	<i>Cute and adorable ferret wizard, wearing coat and suit, steampunk, lantern, anthropomorphic, Jean papististe monge, oil painting</i>	<i>A portrait of a bear wearing a suit in the style of a Baroque painting</i>	<i>Photo of a cat eating a burger like a person</i>	<i>An evil villain holding a mini earth</i>	<i>cinematic still of highly reflective stainless steel train in the desert, at sunset</i>	<i>A cat in a space suit walking on the moon</i>	<i>rural house with a garden and a swimming pool</i>
Stable Diffusion v1.4							
PRDP							
	<i>cubic building on clouds of colorful trees</i>	<i>monkey climbing a skyscraper</i>	<i>cinematic still of an adorable walking robot in the desert, at sunset</i>	<i>Harry potter as a cat, pixar style, octane render, HD, high-detail</i>	<i>cozy house to live in a mountain</i>	<i>a landscape with a river running down the middle in a forest with a sunset behind distant mountains</i>	<i>Horses running on the Great Wall at sunset</i>
Stable Diffusion v1.4							
PRDP							
	<i>A cute blue cat.</i>	<i>a close up of a cat wearing a pikachu hat, reddit, gif, real life charmander, very aesthetic!!!!, soft!!!</i>	<i>Title image: A heartwarming illustration of a cute lion cub named Ladi and a very small koala named Carlo sitting together in the jungle, with an adventurous landscape unfolding in the background. Disney style</i>	<i>wonderful image of a landscape and a medieval tower</i>	<i>futuristic grand fort made out of white marble and extremely intricate carvings across the structure on a martian mountain with fountains and greenery all around</i>	<i>A war weary hamster soldier</i>	<i>An abandoned Segway in the forest</i>
Stable Diffusion v1.4							
PRDP							

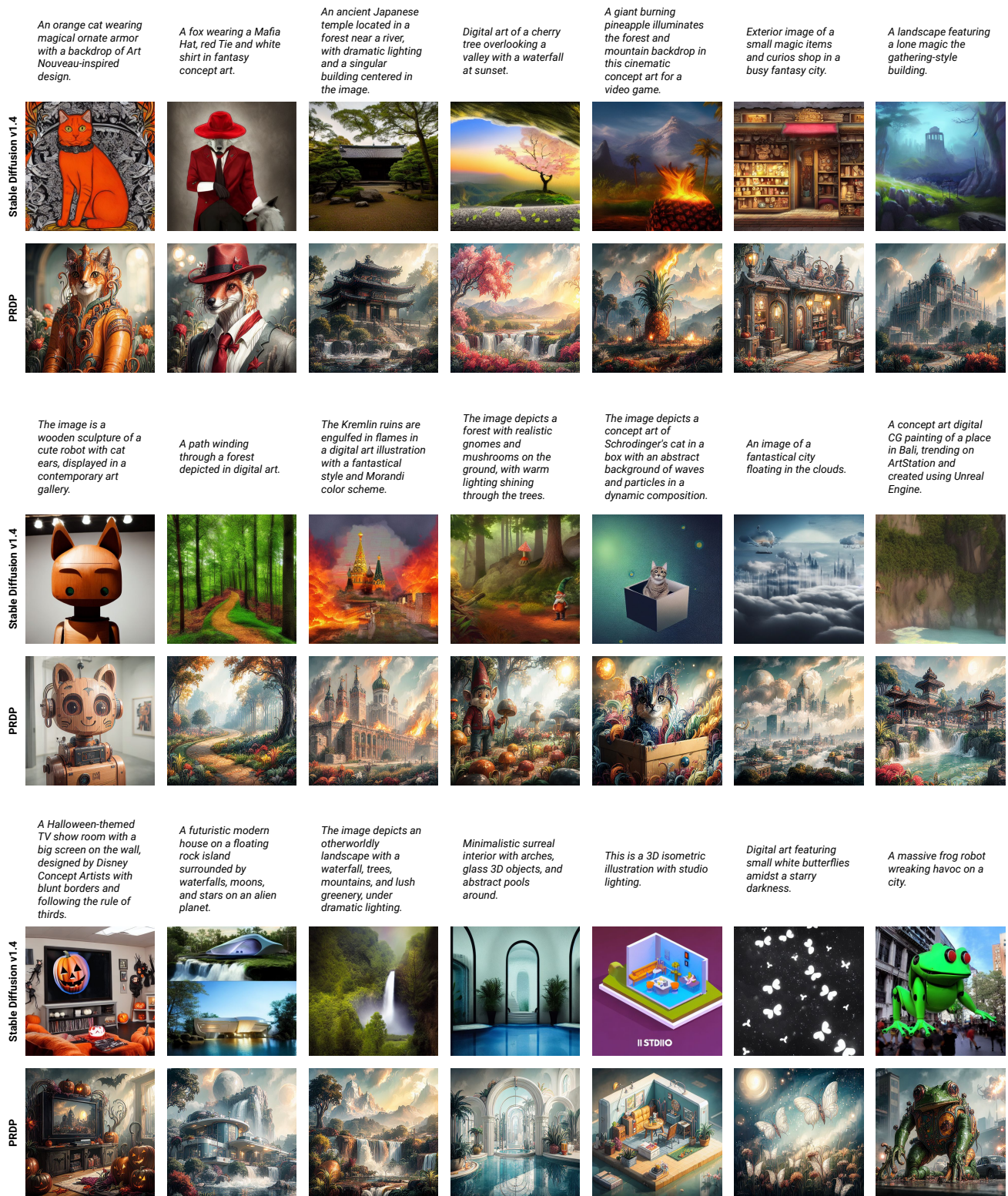
**Figure 11. Generation samples on unseen prompts from the Pick-a-Pic v1 test set.** PRDP is finetuned from Stable Diffusion v1.4 on the training set prompts of Pick-a-Pic v1 dataset, using a weighted combination of rewards: PickScore = 10, HPSv2 = 2, Aesthetic = 0.05. For each prompt, the generation sample from Stable Diffusion v1.4 and PRDP use the same random seed.



	<i>A fox wearing a yellow dress.</i>	<i>A portrait of a silver and white brindle persian cat dressed as a renaissance queen, standing atop a skyscraper overlooking a city.</i>	<i>A cute anthropomorphic fox knight wearing a cape and crown in pale blue armor.</i>	<i>A digital painting of an anthropomorphic corgi lifting weights in a dim gym with intricate details and a dynamic pose.</i>	<i>A toad baby sitting in a rose blossom, depicted in a humorous and detailed illustration.</i>	<i>A chibi frog character surfing at the beach.</i>	<i>An anthropomorphic cat wearing sunglasses and a leather jacket rides a Harley Davidson in Arizona.</i>
Stable Diffusion v1.4							
PRDP							
	<i>Digital art of a female marten animal cartoon character wearing jewelry with a blonde hairstyle.</i>	<i>A bear in an astronaut suit sits on a rock on Mars surrounded by flowers under a starry sky.</i>	<i>A pikachu in a forest illustration.</i>	<i>A portrait of a cat wearing a samurai helmet.</i>	<i>A cute little anthropomorphic bear knight wearing a cape and crown in pale blue armor.</i>	<i>A colorful cartoon tent in a bazaar with a borderlands-inspired aesthetic.</i>	<i>A knitted Capybara wearing sunglasses sips a Mojito at the beach during sunset.</i>
Stable Diffusion v1.4							
PRDP							
	<i>A cartoon satanic priest depicted as an anthropomorphic lamb in a highly detailed 3D render.</i>	<i>The image is a humorous illustration of a furry alien chick nesting in a floral cup.</i>	<i>A landscape with a Maya-style building and Winnie the Pooh on grass.</i>	<i>A fluffy chick is nested in an antique coffee cup in a humorous illustration.</i>	<i>A Fortnite poster featuring chibi kittens wearing cyberpunk headphones and shades, with anime-stylized art by Takeshi Murakami.</i>	<i>A corgi dressed as a bee costume.</i>	<i>A blue bear wearing cowboy boots.</i>
Stable Diffusion v1.4							
PRDP							











































**Figure 12. Generation samples on unseen prompts from the HPD v2 animation benchmark.** PRDP is finetuned from Stable Diffusion v1.4 on the training set prompts of Pick-a-Pic v1 dataset, using a weighted combination of rewards: PickScore = 10, HPSv2 = 2, Aesthetic = 0.05. For each prompt, the generation sample from Stable Diffusion v1.4 and PRDP use the same random seed.





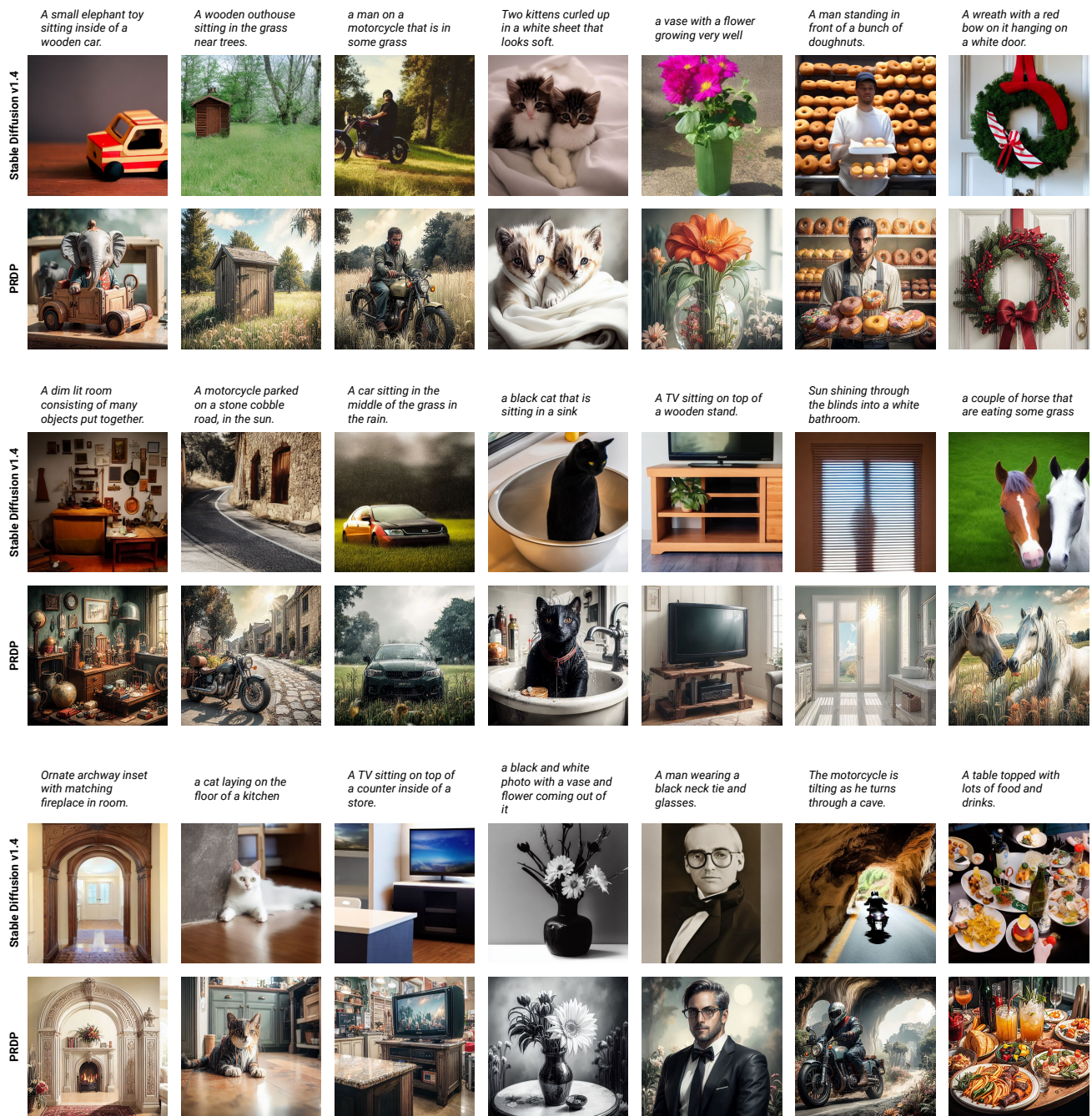
**Figure 13. Generation samples on unseen prompts from the HPD v2 concept art benchmark.** PRDP is finetuned from Stable Diffusion v1.4 on the training set prompts of Pick-a-Pic v1 dataset, using a weighted combination of rewards: PickScore = 10, HPSv2 = 2, Aesthetic = 0.05. For each prompt, the generation sample from Stable Diffusion v1.4 and PRDP use the same random seed.



	<i>A painting of a Persian cat dressed as a Renaissance king, standing on a skyscraper overlooking a city.</i>	<i>A digital painting of a fantasy kitchen environment with elements of cartoons, comics, and manga.</i>	<i>Colorful illustration of a forest tunnel illuminated by sunlight and filled with wildflowers.</i>	<i>A detailed painting of a futuristic spaceship with ornamental features.</i>	<i>The image features a surreal fox and skulls in highly detailed, liquid oilpaint style.</i>	<i>A fluffy owl sits atop a stack of antique books in a detailed and moody illustration.</i>	<i>The image features a castle surrounded by a dreamy garden with roses and a cloudy sky in the background.</i>
Stable Diffusion v1.4							
PRDP							
	<i>A digital painting of a magical ritual location with volumetric lighting and elements from various artworks and games.</i>	<i>An oil painting of a vintage rally car, including a yellow Porsche with smoke and dirt from drifting.</i>	<i>A brownstone building located in a forest setting, painted by Eytan Zana.</i>	<i>A serene meadow with a tree, river, bridge, and mountains in the background under a slightly overcast sunrise sky.</i>	<i>A landscape featuring a unique digital painting-style building.</i>	<i>A night scene of a lavender field with a town and church in the background, reminiscent of Vincent van Gogh's style.</i>	<i>A watercolor painting of a galaxy in a jar.</i>
Stable Diffusion v1.4							
PRDP							
	<i>A landscape with an art nouveau building.</i>	<i>A painting of a girl standing on a mountain looking out at an approaching storm over the ocean, with wind blowing and ocean mist, surrounded by lightning.</i>	<i>A surreal cat with a smile and intricate details.</i>	<i>The image features an ancient Chinese landscape with a mountain, waterfalls, willow trees, and arch bridges set against a blue background.</i>	<i>A digital painting of a blue-skinned wizard with intricate and elegant details, created by multiple artists and posted on Artstation.</i>	<i>A train crosses a trestle bridge in the mountains in an optimistic and vibrant illustration.</i>	<i>A solar eclipse is depicted over a field of grass and flowers with a small forest in the distance, as a matte painting on Art Station by Simon Stalenhag.</i>
Stable Diffusion v1.4							
PRDP							

**Figure 14. Generation samples on unseen prompts from the HPD v2 painting benchmark.** PRDP is finetuned from Stable Diffusion v1.4 on the training set prompts of Pick-a-Pic v1 dataset, using a weighted combination of rewards: PickScore = 10, HPSv2 = 2, Aesthetic = 0.05. For each prompt, the generation sample from Stable Diffusion v1.4 and PRDP use the same random seed.





**Figure 15. Generation samples on unseen prompts from the HPD v2 photo benchmark.** PRDP is finetuned from Stable Diffusion v1.4 on the training set prompts of Pick-a-Pic v1 dataset, using a weighted combination of rewards: PickScore = 10, HPSv2 = 2, Aesthetic = 0.05. For each prompt, the generation sample from Stable Diffusion v1.4 and PRDP use the same random seed.