

UNDERSTANDING SAMPLER STOCHASTICITY IN TRAINING DIFFUSION MODELS FOR RLHF

Anonymous authors

Paper under double-blind review

ABSTRACT

Reinforcement Learning from Human Feedback (RLHF) is increasingly used to fine-tune diffusion models, but a key challenge arises from the mismatch between stochastic samplers used during training and deterministic samplers used during inference. In practice, models are fine-tuned using stochastic SDE samplers to encourage exploration, while inference typically relies on deterministic ODE samplers for efficiency and stability. This discrepancy induces a **reward gap**, raising concerns about whether high-quality outputs can be expected during inference. In this paper, we theoretically characterize this reward gap and provide non-vacuous bounds for general diffusion models, along with sharper convergence rates for Variance Exploding (VE) and Variance Preserving (VP) Gaussian models. Methodologically, we adopt the generalized denoising diffusion implicit models (gDDIM) framework to support arbitrarily high levels of stochasticity, preserving data marginals throughout. Empirically, our findings through large-scale experiments on text-to-image models using denoising diffusion policy optimization (DDPO) and mixed group relative policy optimization (MixGRPO) validate that reward gaps consistently narrow over training, and ODE sampling quality improves when models are updated using higher-stochasticity SDE training.

1 INTRODUCTION

Diffusion models (e.g., Stable Diffusion (Rombach et al., 2022), SDXL (Podell et al., 2024), FLUX (Black Forest Labs, 2024)) have shown strong performance in text-to-image (T2I) tasks, and have also been extended beyond images to video (Ho et al., 2022) and audio (Liu et al., 2023). To meet downstream objectives such as aesthetics, safety, and alignment, it is essential to post-train with RLHF (Ouyang et al., 2022) for preference-driven improvements, often with a KL-regularization term to preserve performance on pretrained tasks (Schulman et al., 2017). Widely used RLHF algorithms include DDPO (Black et al., 2024) and GRPO (Shao et al., 2024) variants (FlowGRPO (Liu et al., 2025), DanceGRPO (Xue et al., 2025), MixGRPO (Li et al., 2025)). DDPO directly optimizes human-preference rewards by casting the denoising process as a Markov Decision Process (MDP); GRPO variants use group-relative advantages. RLHF training may also exhibit unstable trajectories, long inference times, and vulnerability to reward hacking (Skalse et al., 2022; Lee et al., 2025) and efficient, robust samplers (Lu et al., 2022; 2025) are desired for stable, high-quality fine-tuned models. See (Winata et al., 2025) for a broader review of successful RLHF algorithms for generative models.

The classical scheme DDPM (Ho et al., 2020), a discretization of the score-based backward SDE (Risken, 1996; Song et al., 2021b), allows for trajectory diversity and rich exploration in RLHF training; the deterministic DDIM sampler (Song et al., 2021a) follows a probability-flow ODE and facilitates inference. High stochasticity during training is crucial for effective exploration of the reward landscape. However, during inference, deterministic ODE sampling is preferred to ensure consistency and computational efficiency. This creates an inherent tension: we train via a highly stochastic process (SDE) but we deploy a deterministic process (ODE) for inference. This discrepancy necessitates the following question:

To what extent can we guarantee ODE inference quality after we fine-tune diffusion models with SDE sampled trajectories?

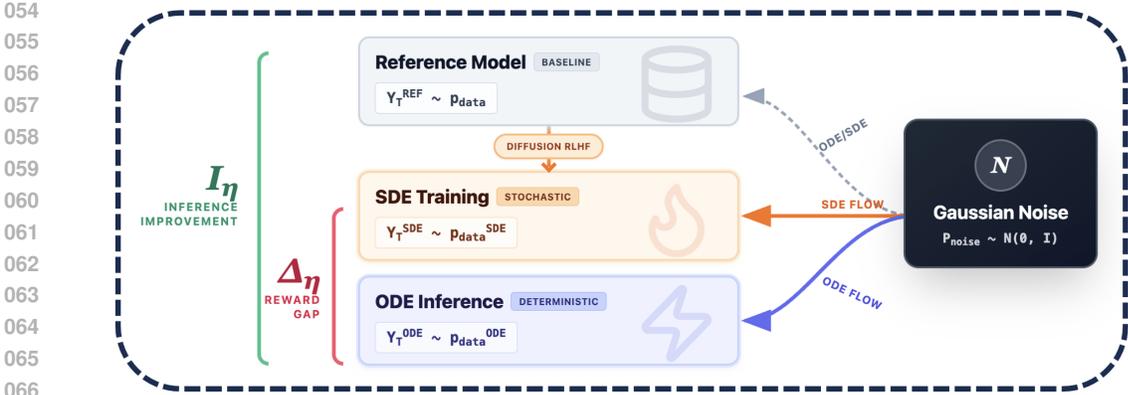


Figure 1: An illustration of the problem formulation. Δ_η represents the “Reward Gap” and I_η represents the ODE reward improvement after SDE fine-tune. We would like $\Delta_\eta \downarrow$ with $I_\eta \uparrow$.

Figure 1 illustrates the problem formulation. While SDE and ODE denoising processes (equation 3) yield the same terminal marginal for the pretrained model (see Section 2.1), a *distributional difference* arises between the distribution ($p_{\text{data}}^{\text{SDE}}$) after stochastic fine-tuning and the deterministic inference outcome ($p_{\text{data}}^{\text{ODE}}$) (see Section 3). We denote η as the stochasticity scale, where $\eta = 1.0$ corresponds to the standard variance in DDPM. We define this discrepancy as the **Reward Gap** (Δ_η), which must be bounded to preserve the ODE reward improvement (I_η) gained over the baseline. In this paper, we present our analysis from three perspectives:

- **Theory:** We bound Δ_η between a generally SDE-fine-tuned model and its ODE-sampling counterpart using Gronwall’s inequality. Specifically, for VE, VP and mixture Gaussian models, the gap shrinks at sharp rates $O(1/T)$ and $O(e^{-T^2}/2)$, where T is the denoising time horizon.
- **Methodology:** To handle high stochasticity beyond $\eta = 1.0$ without changing marginals, we adopt the gDDIM scheme (Zhang et al., 2023) for arbitrary stochasticity levels. This framework is introduced in Section 2.2 and 5. See Appendix E for more technical details.
- **Experiments:** We evaluate large-scale T2I models and RLHF algorithms, DDPO and Mix-GRPO, to quantify the reward gap across multiple reward functions. The results show that Δ_η decrease as training quality improves and $\eta = 1.2$ yields superior ODE inference performance, I_η .

Related Literature. Higher-order ODE solvers such as RX-DPM (Choi et al., 2025), DEIS (Zhang & Chen, 2022) approximate probability-flow. Recent analyses give broader convergence guarantees (Huang et al., 2025) and better representability (Chen et al., 2023b) for score-based diffusion. In parallel, RL-based fine-tuning leverages stochastic sampling: DRaFT differentiates through noisy trajectories (Clark et al., 2023), Score-as-Action casts fine-tuning as stochastic control (Zhao et al., 2025), and Adjoint Matching enforces optimal memoryless noise schedules (Domingo-Enrich et al., 2024). Large-scale studies combine multiple rewards (Zhang et al., 2024a); SEPO, D3PO, and ImageReFL boosts alignment (Zekri & Boullé, 2025; Yang et al., 2024; Sorokin et al., 2025). Finally, (Liang et al., 2025) gives discretization error bound and (Wu et al., 2024) discusses guidance (Ho & Salimans, 2021) for Gaussian Mixture priors.

Organization of the paper. The remainder of the paper is organized as follows. In Section 2, we provide background on diffusion RLHF and gDDIM in particular. Section 3 computes *exact* Δ_η and I_η for simple controls. More complex networks are analyzed in Section 4. Numerical experiments in Section 5 validate our theoretical bounds by fine-tuning large-scale T2I models with DDPO and GRPO algorithms. Conclusion is given in Section 6.

2 PRELIMINARIES

2.1 FORWARD-BACKWARD DIFFUSION MODELS

The goal of diffusion models (Song et al., 2021b) is to generate new samples similar to $p_{\text{data}}(\cdot)$. The forward process is given by an SDE:

$$dX_t = f(t, X_t)dt + g(t)dB_t, \quad X_0 \sim p_{\text{data}}(\cdot), \quad (1)$$

where $\{B_t\}$ is the d -dimensional Brownian motion, and $f : \mathbb{R}_+ \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ are given model parameters under regularity conditions (Stroock & Varadhan, 1979).

We use *denoising score matching* (Vincent, 2011) to minimize the mean square error between a parametrized neural network $s_\theta(t, X_t)$ and the unknown term $\nabla \log p(T-t, \bar{X}_t)$ in the time reversal SDE (Haussmann & Pardoux, 1986). The reverse-time SDE becomes:

$$dY_t = (-f(T-t, Y_t) + g^2(T-t)s_\theta(T-t, Y_t))dt + g(T-t)dB_t, \quad Y_0 \sim p_{\text{noise}}(\cdot). \quad (2)$$

For richer levels of *sampler stochasticity*, a flexible stochasticity parameter $\eta \geq 0$ is introduced:

$$dY_t = \left(-f(T-t, Y_t) + \frac{1+\eta^2}{2}g^2(T-t)s_\theta(T-t, Y_t) \right) dt + \eta g(T-t)dB_t, \quad Y_0 \sim p_{\text{noise}}(\cdot), \quad (3)$$

With a divergence-free noise scale, equation 3 preserves the terminal marginals (Anderson, 1982).

2.2 DISCRETIZATION – STOCHASTIC GDDIM

As shown in (Song et al., 2021b), for most diffusion models, the model parameters are:

$$f(t, x) = \frac{1}{2} \frac{d \log \alpha_t}{dt} x \quad \text{and} \quad g(t) = \sqrt{-\frac{d \log \alpha_t}{dt}},$$

where $\{\alpha_t\}_{t=0}^T$ is a decreasing sequence from 1 to 0. For all $\eta \geq 0$, (Zhang et al., 2023) proposed the *generalized DDIM* update:

$$x_{t-\Delta t} = \sqrt{\frac{\alpha_{t-\Delta t}}{\alpha_t}} x_t + \left(\sqrt{\frac{\alpha_{t-\Delta t}}{\alpha_t}}(1-\alpha_t) - \sqrt{(1-\alpha_{t-\Delta t}-\sigma_t^2)(1-\alpha_t)} \right) s_\theta(t, x_t) + \sigma_t(\eta) \mathcal{N}(0, I), \quad x_T \sim p_{\text{noise}}(\cdot). \quad (4)$$

where $\{x_t\}_{t=0}^T$ is the discretized sequence of backward diffusion Y_t and

$$\sigma_t(\eta) = (1-\alpha_{t-\Delta t}) \left(1 - \left(\frac{1-\alpha_{t-\Delta t}}{1-\alpha_t} \right)^{\eta^2} \left(\frac{\alpha_t}{\alpha_{t-\Delta t}} \right)^{\eta^2} \right). \quad (5)$$

For any $\eta > 0$, the discrete scheme equation 4 has terminal marginal equal to equation 3. We use this discretization formulation for post-training with $\eta > 1.0$ in Section 5. Details in Appendix E.

2.3 DIFFUSION RLHF

RLHF was originally proposed for LLM alignment (Bai et al., 2022; Ouyang et al., 2022). (Black et al., 2024; Fan et al., 2023; Lee et al., 2023) first formulated the denoising step as Markov decision processes (MDPs) to fine-tune diffusions. (Gao et al., 2024; Zhao et al., 2024; 2025) developed an RL approach based on continuous-time models. Here we focus on *Denoising Diffusion Policy Optimization* (DDPO) and *Group Relative Policy Optimization* (GRPO).

DDPO (Black et al., 2024): We formulate the denoising steps $\{x_T, x_{T-1}, \dots, x_0\}$ as an MDP:

$$\begin{aligned} s_t &:= (\mathbf{c}, t, x_t), & \pi(a_t | s_t) &:= p_\theta(x_{t-1} | x_t, \mathbf{c}), & \mathbb{P}(s_{t+1} | s_t, a_t) &:= (\delta_{\mathbf{c}}, \delta_{t-1}, \delta_{x_{t-1}}), \\ a_t &:= x_{t-1}, & \rho_0(s_0) &:= (p(\mathbf{c}), \delta_T, \mathcal{N}(0, I)), & R(s_t, a_t) &:= \begin{cases} r(x_0, \mathbf{c}) & \text{if } t = 0, \\ 0 & \text{otherwise,} \end{cases} \end{aligned}$$

where δ_\bullet is the Dirac point mass concentrating at \bullet . The objective of DDPO is to maximize:

$$\mathcal{J}_{DDPO}(\theta) := \mathbb{E}_{\mathbf{c} \sim p(\mathbf{c}), x_0 \sim p_\theta(x_0|\mathbf{c})} [r(x_0, \mathbf{c})]. \quad (6)$$

The policy gradient of equation 6 is $\nabla_\theta \mathcal{J}_{DDPO} = \mathbb{E} \left[\sum_{t=0}^T \nabla_\theta \log p_\theta(x_{t-1}|x_t, \mathbf{c}) r(x_0, \mathbf{c}) \right]$ (Williams, 1992). A common parameterization of p_θ is isotropic Gaussian (Mohamed et al., 2020).

GRPO (Shao et al., 2024): GRPO first computes the *advantages*:

$$A_i := \frac{r(x_0^i, \mathbf{c}) - \frac{1}{G} \sum_{k=1}^G r(x_0^k, \mathbf{c})}{\text{std}(\{r(x_0^k, \mathbf{c})\}_{k=1}^G)},$$

where G is the group size; $\text{std}(\{r(x_0^k, \mathbf{c})\}_{k=1}^G)$ denotes the *standard deviation* of group rewards. The objective of GRPO is to maximize:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{N} \sum_{n=1}^N \min(\rho_t^i(\theta) A_i, \text{clip}(\rho_t^i(\theta), 1 - \epsilon, 1 + \epsilon) A_i) \right], \quad (7)$$

where $\text{clip}(\rho_t^i(\theta), 1 - \epsilon, 1 + \epsilon) := \min\{\max\{\rho_t^i(\theta), 1 - \epsilon\}, 1 + \epsilon\}$ restricts the likelihood ratio $\rho_t^i(\theta) = \frac{p_\theta(x_{t-1}^i|x_t^i, \mathbf{c})}{p_{\theta_{\text{old}}}(x_{t-1}^i|x_t^i, \mathbf{c})}$ within the range $[1 - \epsilon, 1 + \epsilon]$ by imposing hard constraints on the boundaries; $\{x_t^i\}$ is the i -th sample trajectory in the group; N is the number of grouped outputs, and the expectation is with respect to $\mathbf{c} \sim p(\mathbf{c})$ and $\{x_t^i\} \sim \pi_{\theta_{\text{old}}}(\cdot|\mathbf{c})$.

3 THEORY ON THE REWARD GAP

As discussed in Section 1, forward-backward diffusion models typically employ *stochastic* SDE dynamics for exploration in the RLHF training, while adopting *deterministic* ODE samplers for inference (see Figure 1). In practice, we follow entropy-regularized objectives to fine-tune a parametrized family of score functions (Uehara et al., 2024; Tang, 2024):

Definition 3.1. *Define*

1. $\{Y_T^{\text{REF}}\}$ as the samples extracted from the referenced base model following equation 3 with terminal marginal p_{data} ;
2. We fine-tune the following objectives under a fixed stochasticity parameter η :

$$F_\eta(\theta) = \mathbb{E}[r(Y_T^\theta)] - \beta \text{KL}(Y_T^\theta \| Y_T^{\text{REF}}), \quad \theta_\eta^* := \arg \max_\theta F_\eta(\theta); \quad (8)$$

let $\{Y_t^{\text{SDE}}\} = \{Y_t(\theta_\eta^*)\}$ be the optimal fine-tuned model, and $\{Y_t^{\text{ODE}}\}$ be the deterministic sampler obtained by letting $\eta = 0$:

$$\begin{aligned} dY_t^{\text{SDE}} &= \left(-f(T-t, Y_t^{\text{SDE}}) + \frac{1+\eta^2}{2} g^2(T-t) s_{\theta_\eta^*}(T-t, Y_t^{\text{SDE}}) \right) dt + \eta g(T-t) dB_t, \\ dY_t^{\text{ODE}} &= \left(-f(T-t, Y_t^{\text{ODE}}) + \frac{1}{2} g^2(T-t) s_{\theta_\eta^*}(T-t, Y_t^{\text{ODE}}) \right) dt. \end{aligned} \quad (9)$$

Although a pretrained model admits identical terminal marginals for arbitrary η in the form of equation 3, once the data distribution is aligned to a human-preference reward $r(x)$, the distribution of $\{Y_T^{\text{ODE}}\}$ and $\{Y_T^{\text{SDE}}\}$ (we denote as $p_{\text{data}}^{\text{SDE}}$ and $p_{\text{data}}^{\text{ODE}}$) are not necessarily the same.

Under a fixed noise level η , the terminal marginal satisfies a *reward-tilting* distribution $p_{\text{data}}^{\text{SDE}} \propto p_{\text{data}}(x) \exp(r(x)/\beta)$ (Zhang et al., 2024b). The terminal score function

$$s_{\theta_\eta^*}(T, x) \approx \nabla \log p_{\theta_\eta^*}(T, x) = \nabla \log p_{\text{data}}(x) + \frac{\nabla r(x)}{\beta},$$

but the fine-tuned parameter θ_η^* does not regulate the intermediate trajectory, as the RLHF objective equation 8 only considers the terminal marginal. Therefore, we could not establish a general relationship between $s_{\theta_\eta^*}(t, x)$, $r(x)$, and $\nabla \log p_{\theta_\eta^*}(t, x)$ for $0 < t < T$. In general,

$$\nabla \cdot (p_{\theta_\eta^*}(t, x) \cdot s_{\theta_\eta^*}(t, x)) \neq \Delta p_{\theta_\eta^*}(t, x).$$

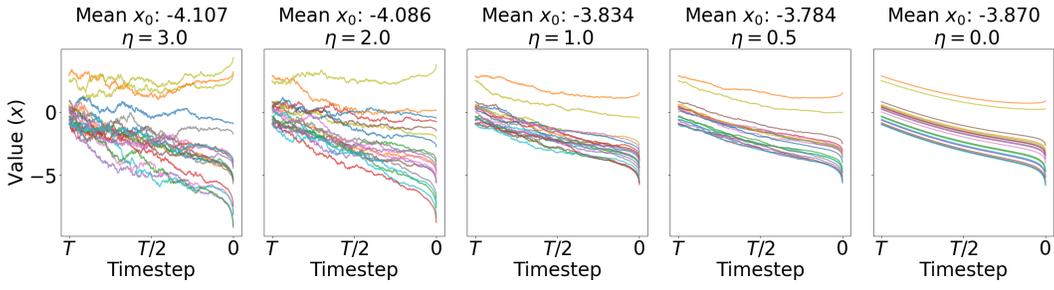


Figure 2: One-dimensional VP dynamics under different noise levels η , using the same control function, demonstrating $\mathbb{E}(Y_T^{\text{SDE}}) \approx \mathbb{E}(Y_T^{\text{ODE}})$ despite $p_{\text{data}}^{\text{SDE}} \neq p_{\text{data}}^{\text{ODE}}$.

As the reverse-time Fokker–Planck correspondence no longer satisfies the divergence free condition (Anderson, 1982), $p_{\text{data}}^{\text{SDE}}$ and $p_{\text{data}}^{\text{ODE}}$ need not coincide.

However, we are motivated by an observation in Figure 2 that 1D Variance Preserving (VP) dynamics, when run with the same constant downward drift, maintain a similar mean value regardless of the noise level η . We next rigorously quantify and generalize this observation. Bounded reward gaps for simple dynamics are important for two reasons:

- **Error Isolation.** Under VP (as well as VE) models, we are able to compute the *exact* reward difference induced by changing η , free from the well-studied score-approximation and discretization error (Liang et al., 2025), with which we will also analyze in Section 4 and upper-bound reward gaps in those more complex settings.
- **Parametrization Insight.** Parametrization is inherent to the pre-trained (see equation 2) and the fine-tuned model (see equation 9) and affect RLHF performances in practice (Han et al., 2025). Here we show that reward gap is small *even under simple parametrization*. As Section 4 bounds complex controls *quadratic* in η , our results are valid for general networks.

Our goals are to quantify the reward improvement of fine-tuning I_η and the reward gap Δ_η between SDE and ODE inference. Here we formally define:

Definition 3.2. As illustrated in Figure 1, we define

1. *Improvement of fine-tuning as the reward difference between $\{Y_t^{\text{ODE}}\}$ and $\{Y_t^{\text{REF}}\}$,* $I_\eta := \mathbb{E}[r(Y_T^{\text{ODE}})] - \mathbb{E}[r(Y_T^{\text{REF}})] \geq 0$.
2. *Reward gap as the reward difference between $\{Y_t^{\text{ODE}}\}$ and $\{Y_t^{\text{SDE}}\}$,* $\Delta_\eta := |\mathbb{E}[r(Y_T^{\text{ODE}})] - \mathbb{E}[r(Y_T^{\text{SDE}})]|$.

In what follows, we analyze VE/VP models with a Gaussian or mixture Gaussian target distribution, where the score function has a closed-form expression without approximation and discretization errors, so the only discrepancy comes from changing η .

Remark. $\mathbb{E}[r(Y_T^{\text{ODE}})] \geq \mathbb{E}[r(Y_T^{\text{REF}})]$ holds for all models, following from the optimality of θ_η^* . However, the sign of $(\mathbb{E}[r(Y_T^{\text{ODE}})] - \mathbb{E}[r(Y_T^{\text{SDE}})])$ depends on the process dynamics as well as RLHF algorithms.

3.1 VE WITH A GAUSSIAN TARGET

We first consider the one-dimensional VE model:

$$dX_t = \sqrt{2t} dB_t, \quad \text{with } p_{\text{data}}(\cdot) = \mathcal{N}(0, 1). \quad (10)$$

Since $X_t \sim \mathcal{N}(0, t^2 + 1)$, the (exact) score function is $\nabla \log p(t, x) = -\frac{x}{t^2 + 1}$. So the “pretrained” model is:

$$dY_t^{\text{REF}} = -\frac{(1 + \eta^2)(T - t)}{1 + (T - t)^2} Y_t^{\text{REF}} dt + \eta \sqrt{2(T - t)} dB_t.$$

Next we set the reward function $r(x) = -(x - 1)^2$, so the goal of fine-tuning is to drive the sample towards mean 1. We also specify the fine-tuned SDE and ODE as defined in equation 9 by

$$\begin{aligned} dY_t^{\text{SDE}} &= -\frac{(1 + \eta^2)(T - t)}{1 + (T - t)^2}(Y_t^{\text{SDE}} + \theta_\eta^*)dt + \eta\sqrt{2(T - t)}dB_t, \\ dY_t^{\text{ODE}} &= -\frac{T - t}{1 + (T - t)^2}(Y_t^{\text{ODE}} + \theta_\eta^*)dt, \end{aligned}$$

with θ_η^* maximizing the entropy-regularized reward (equation 8). The following theorem gives bounds for I_η and Δ_η under VE, and the proof is deferred to Appendix C.1.

Theorem 3.1. Consider the Variance Exploding model (equation 10), with the reward function

$$r(x) = -(x - 1)^2. \quad (11)$$

For $\eta > 0$, we have $\theta_\eta^* = -\left[\left(1 + \frac{\beta}{2}\right)\left(1 - (1 + T^2)^{-\frac{1+\eta^2}{2}}\right)\right]^{-1}$. Moreover, for $T \geq 1$

$$0 \leq \Delta_\eta \leq \frac{1}{2T} + o\left(\frac{1}{T}\right) \quad \text{and} \quad \mathcal{I}_\eta \geq 1 - \frac{1}{2T} + o\left(\frac{1}{T}\right). \quad (12)$$

3.2 VP WITH A GAUSSIAN TARGET

Now we consider the one-dimensional VP model:

$$dX_t = -tX_tdt + \sqrt{2t}dB_t, \quad \text{with } p_{\text{data}}(\cdot) = \mathcal{N}(0, 1). \quad (13)$$

Under the same setup as equation 9 and the VE case, we have:

$$\begin{aligned} dY_t^{\text{REF}} &= -\eta^2(T - t)Y_t^{\text{REF}}dt + \eta\sqrt{2(T - t)}dB_t, \\ dY_t^{\text{SDE}} &= -\eta^2(T - t)Y_t^{\text{SDE}}dt - (1 + \eta^2)(T - t)\theta_\eta^*(t)dt + \eta\sqrt{2(T - t)}dB_t, \\ dY_t^{\text{ODE}} &= -(T - t)\theta_\eta^*(t)dt, \end{aligned}$$

where a time-dependent control $\theta_\eta(t) := \theta_\eta e^{-\frac{(T-t)^2}{2}}$ is used for fine-tuning. The following theorem gives bounds for I_η and Δ_η under VP, and the proof is deferred to Appendix C.2.

Theorem 3.2. Consider the VP model (equation 13), with the reward function $r(x) = -(x - 1)^2$.

For $\eta > 0$, we have $\theta_\eta^* = -\left[\left(1 + \frac{\beta}{2}\right)\left(1 - e^{-\frac{(1+\eta^2)T^2}{2}}\right)\right]^{-1}$. Moreover, for $T \geq 1$

$$0 \leq \Delta_\eta \leq \frac{e^{-T^2}}{2} + o\left(e^{-T^2}\right) \quad \text{and} \quad \mathcal{I}_\eta \geq 1 - \frac{e^{-T^2}}{2} + o\left(e^{-T^2}\right). \quad (14)$$

3.3 VE/VP WITH A MIXTURE GAUSSIAN TARGET

The previous results can be extended to multidimensional setting, with a mixture Gaussian target distribution. Recall that the probability density of a mixture Gaussian has the form:

$$\sum_{i=1}^k \frac{\alpha_i}{(2\pi)^{d/2}(\det \Sigma_i)^{1/2}} \cdot \exp\left(-\frac{1}{2}(x - \mu_i)^T \Sigma_i (x - \mu_i)\right),$$

where α_i is the weight of the i -th Gaussian component. The following corollary bounds the reward gap for a mixture Gaussian target distribution, and the proof is deferred to Appendix C.3.

Corollary 3.1. Let the reward function be $r(x) = -\|x - \mathbf{r}\|^2$ such that $\mu_i \cdot \mathbf{r} = 0$ for all $i \in \{1, \dots, k\}$, $\Sigma_i \equiv \mathbf{I}_d$ and $\mathbb{E}[Y_T^{\text{REF}}] = \mathbf{0}$. Then the same bounds on reward gap hold, i.e.,

$$\Delta_\eta \leq \begin{cases} (2T)^{-1} & \text{for VE,} \\ e^{-T^2}/2 & \text{for VP.} \end{cases} \quad (15)$$

For all these examples, the reward gap $\Delta_\eta \downarrow 0$ (independent of η), as $T \rightarrow \infty$. Such a phenomenon will also be observed in fine-tuning T2I models with more complex rewards.

	η	ImageReward	PickScore	HPS_v2	Aesthetic
Base	–	0.32	20.69	0.253	5.38
ImageReward	1.0	0.84 (0.92)	20.89 (20.94)	0.268 (0.271)	5.64 (5.72)
	1.2	0.92 (1.03)	20.95 (20.99)	0.268 (0.289)	5.70 (5.80)
	1.5	0.77 (0.91)	20.90 (20.95)	0.266 (0.269)	5.61 (5.74)
PickScore	1.0	0.59 (0.73)	21.11 (21.28)	0.265 (0.267)	5.56 (5.68)
	1.2	0.59 (0.69)	21.17 (21.33)	0.264 (0.264)	5.57 (5.69)
	1.5	0.57 (0.70)	21.10 (21.36)	0.264 (0.262)	5.60 (5.77)

Table 1: Performance for ODE (SDE) samplers under DDPO fine-tuning. Bold numbers indicate the highest evaluations among all stochasticity, demonstrating $\eta = 1.2$ in general performs the best.

4 NON-VACUOUS BOUND ON THE REWARD GAP

In Section 3, our analysis relies on the explicit computation of score functions with simple controls, which is not available for fine-tuning in general. Here our goal is to bound the Wasserstein-2 distance between $p_{\text{data}}^{\text{SDE}}$ and $p_{\text{data}}^{\text{ODE}}$ with mild assumptions before applying it to a L_r -Lipschitz reward $r(\cdot)$, i.e. $|r(y_1) - r(y_2)| \leq L_r \cdot \|y_1 - y_2\|$.

Assumption 4.1. *The following conditions hold for all y, y_1, y_2 along inference trajectories,*

1. *Lipschitz of s_θ : There exists $L > 0$ such that $\|s_\theta(t, y_1) - s_\theta(t, y_2)\| \leq L\|y_1 - y_2\|$.*
2. *L^2 bound on s_θ : There exists $A > 0$ such that $\sup_{0 \leq t \leq T} \mathbb{E}[\|s_\theta(t, y)\|^2] \leq A$.*

Condition 1 (Lipschitz of scores) and 2 (L^2 bound) are standard in stability analysis. Also, there always exists a continuous-time VP/VE SDE on $[0, 1]$ whose discretization (Song et al., 2021a) matches the large-scale T2I models’ (i.e. Stable Diffusion and SDXL) training schedule. In addition, models (FLUX) based on rectified flow inherently uses a normalized time interval. Therefore, for large η , the following theorem yields $O(\eta^2)$ bound on W_2 distance, see proof in Appendix D.1.

Theorem 4.1. *If Assumption 4.1 hold with $T = 1$, $W_2(Y_T^{\text{ODE}}, Y_T^{\text{SDE}}) \leq C \cdot \max\{\eta, \eta^2\} \cdot \sqrt{1 + A}$, where C only depends on $\|g\|_\infty$ and L .*

Proposition 4.1. *Let Assumption 4.1 hold with L_r -Lipschitz rewards r , $\Delta_\eta \leq L_r \cdot C \cdot \max\{\eta, \eta^2\} \cdot \sqrt{1 + A}$, where C only depends on $\|g\|_\infty$ and L .*

An improved $O(\eta)$ bound relies on contractive coefficients (see Appendix D.2). This contractivity can arise from a contractive drift f (Tang & Zhao, 2024). Alternatively, even though the backward dynamics of classical VP models are expansive, a strongly log-concave terminal distribution (Gao et al., 2025) still satisfies the contractive condition (see Assumption D.1), which can hold if the conditional terminal distribution for a fixed prompt c is approximately unimodal.

Quadratic and linear growth rate on η is consistent with our empirical observation: T2I fine-tuning preserves quality for ODE inference; rewards only deteriorate under very large η (see Table 3).

Remark. Discretization errors can be incorporated into the W_2 bound as an additional term $\epsilon_d(h)$ depending on the time-step size h (Liang et al., 2025), with $\epsilon_d(h) \rightarrow 0$ as $h \rightarrow 0$. Under Assumption 4.1, Euler-Maruyama yields such a vanishing term; see Appendix D.3 for details.

5 NUMERICAL EXPERIMENTS

In this section, we fine-tune large-scale T2I models under large η and examine its reward gap Δ_η with RLHF algorithms DDPO and MixGRPO (see Section 2.3). Preference rewards for fine-tuning DDPO include the LAION aesthetic (Schuhmann et al., 2022), HPS_v2.1 (Wu et al., 2023), PickScore (Kirstain et al., 2023), and ImageReward (Xu et al., 2023). The rewards for fine-tuning MixGRPO include HPS CLIP, PickScore, ImageReward, and Unified Reward (Wang et al., 2025).

According to equation 4 and equation 5, we use a high stochasticity $\eta \geq 1.0$ under **gDDIM scheme** (see Section 2.2 and Appendix E) to generate robust training samples $\{Y_T^{\text{SDE}}\}$ and compare them with deterministic inference samples $\{Y_T^{\text{ODE}}\}$ following equation 9. The stochasticity scale η controls the noise level of backward dynamics.

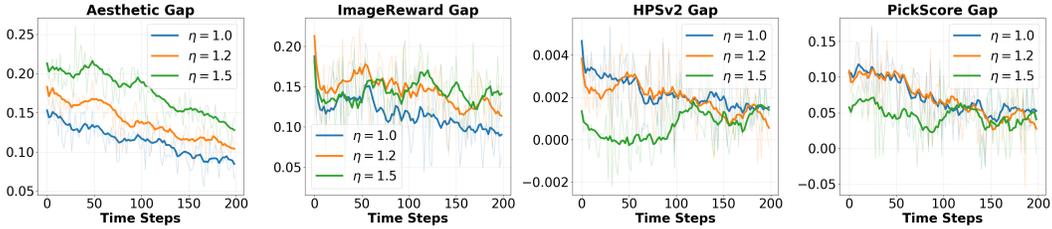


Figure 3: Evolution of reward gap (the reward difference between SDE sampling and ODE sampling) during DDPO training under PickScore fine-tuning with stochasticity $\eta \in \{1.0, 1.2, 1.5\}$. The gap for multiple rewards remains small as training step N progresses.

Remark. In each experiment, the training stochasticity η and the number of training steps N are always specified. Choices of other hyperparameters are detailed in Appendix F.

5.1 DDPO

We use Stable Diffusion v1.5 (Rombach et al., 2022) as the base model, and fine-tune it with DDPO. During training, we adopt ImageReward and PickScore as preference rewards, while Aesthetic and HPS_v2 are included as additional evaluation metrics. Figure 7 (and Appendix H) provides representative generations from the fine-tuned models. Our main observations are summarized as follows:

- **High Stochasticity Benefits Moderate Training Steps.** We compare fine-tuning under ImageReward and PickScore at $\eta \in \{1.0, 1.2, 1.5\}$ for $N = 200$ steps. As shown in Table 1, $\eta = 1.2$ under ImageReward achieves the best in-domain and out-of-domain performance, while PickScore’s performances depend on evaluation metrics.
- **Decreasing Reward Gap with Quality Improvement.** To study the reward gaps for smaller or larger N , we experiment under PickScore with $\eta = 1.2$ and calculate SDE–ODE reward differences (here we report the SDE over ODE performance) under multiple preference functions every 200 steps until reward collapse. As shown in Table 2, the gap decreases as image quality improves for both samplers.
- **Richer Prompt Contents Reduce Reward Gap.** We compare performances with animal versus more comprehensive prompts (see Appendix H.4) under ImageReward with $\eta = 1.2$. As shown in Table 3, complex prompts generate higher post-tuning rewards and higher in-group variance in post-training. Moreover, their richer instructions reduce the SDE–ODE reward gap.

	$N = 0$	200	400	600	800	1000	1200	1400	(1476)
ImgRwrld Gap	0.160	0.119	0.102	0.028	0.057	0.027	0.006	0.020	(0.564)
HPSv2 Gap	0.0047	0.0032	0.0032	0.0018	0.0027	0.0015	-0.0028	0.0011	(0.0308)
Aesthetic Gap	0.162	0.113	0.078	0.077	0.072	0.017	0.030	-0.010	(0.685)
PickScore Gap	0.115	0.146	0.167	0.118	0.178	0.106	0.129	0.090	(0.907)
SDE Reward	20.82	21.31	21.65	21.90	22.07	22.27	22.31	22.42	(18.92)
ODE Reward	20.73	21.17	21.50	21.78	21.88	22.17	22.20	22.32	(17.04)

Table 2: PickScore training with $\eta = 1.2$ until reward collapses. Smallest reward gaps and best sampler performances locate at the large training steps $N = 1200, 1400$.

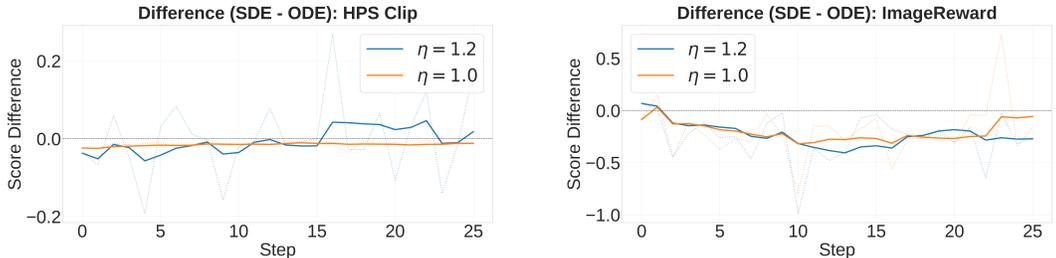
	Animal Prompts			Comprehensive Prompts		
	$N = 0$	100	200	$N = 0$	100	200
Mean	0.38 (0.540)	0.67 (0.81)	0.92 (1.03)	0.35 (0.47)	0.76 (0.87)	1.09 (1.17)
Std	0.81 (0.80)	0.78 (0.72)	0.71 (0.65)	1.03 (1.00)	0.91 (0.87)	0.80 (0.73)
$\Delta_{\eta=1.2}$	0.15	0.14	0.12	0.12	0.11	0.11

Table 3: Performance Comparison between prompts of different complexity with $\eta = 1.2$. More complicated prompts yields faster fine-tuning improvements, larger in-group variances, and smaller SDE–ODE reward gaps.

432 5.2 MixGRPO
433

434 We use FLUX.1 (Black Forest Labs, 2024) as the base model, and fine-tune it with MixGRPO,
435 which is a sliding-window sampler that alternates between ODE and SDE schemes for 25 training
436 steps in total. Training is carried out with multiple rewards combined using equal weights, while
437 evaluation is reported on ImageReward and HPSClip.

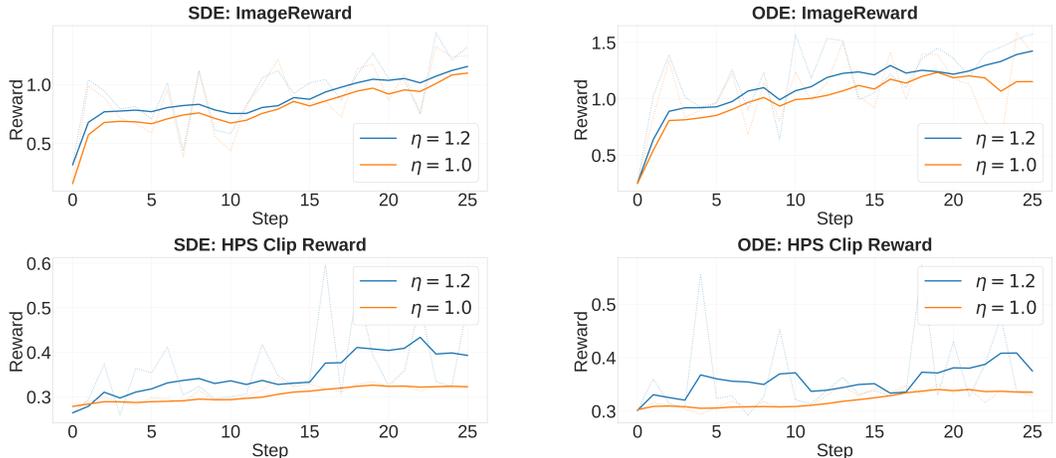
438 • **Bounded Reward Gap Under High Stochasticity.** As shown in Figure 4, Δ_η converges to zero
439 for both HPS Clip and ImageReward. Also, with $\eta = 1.2$ and ImageReward metric, ODE sampling
440 in MixGRPO consistently outperforms the mixed SDE–ODE scheme, in contrast to the results from
441 DDPO.



444

445 **Figure 4:** Bounded reward gap for MixGRPO, displayed with 7-step moving average. Here a positive score
446 difference means the SDE component achieves higher reward than the ODE component.
447

448 • **T2I Quality Improvement.** The upper (evaluated with ImageReward) and lower panels (evaluated
449 with HPS Clip) of Figure 5 further show that samplers trained with larger stochasticity ($\eta = 1.2$)
450 perform consistently better than standard DDIM stochasticity ($\eta = 1.0$) on training prompts. For
451 example, in Figure 6, the generation with $\eta = 1.2$ correctly aligns with the “trapped inside” prompt
452 instruction, whereas the generation with $\eta = 1.0$ fails to do so.
453



456

457 **Figure 5:** Performance improvement for MixGRPO, displayed with 7-step moving average.
458

459

478 6 CONCLUSION AND FURTHER DIRECTIONS
479

480 This work clarifies the tension between stochastic SDE training and deterministic ODE inference in
481 diffusion RLHF. By proving a bounded reward gap Δ_η and empirically showing that higher training
482 stochasticity (e.g., $\eta = 1.2$) improves deterministic image quality, we provide theoretical support
483 for “training with SDE, inference with ODE”. Future work includes quantifying how distribution
484 shift and the choice of reward function separately affect reward gaps.
485



Figure 6: Comparison of ODE image generation by FLUX with MixGRPO fine-tuning, stochasticity $\eta = 1.2$ (below) and $\eta = 1.0$ (above). Higher stochasticity shows better alignments to details.

Prompts (from left to right): “A steampunk pocketwatch owl is trapped inside a glass jar buried in sand, surrounded by an hourglass and swirling mist.”, “An androgynous glam rocker poses outside CBGB in the style of Phil Hale.”, “A digital painting by Loish featuring a rush of half-body, cyberpunk androids and cyborgs adorned with intricate jewelry and colorful holographic dreads.”

REFERENCES

- Brian D. O. Anderson. Reverse-time diffusion equation models. *Stochastic Process. Appl.*, 12(3): 313–326, 1982.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, and et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. In *ICLR*, 2024.
- Black Forest Labs. Flux. *GitHub repository*, 2024. URL <https://github.com/black-forest-labs/flux>.
- Hongrui Chen, Holden Lee, and Jianfeng Lu. Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *Proceedings of the 40th International Conference on Machine Learning*, 2023a.
- Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R. Zhang. Sampling is as easy as learning the score: Theory for diffusion models with minimal data assumptions. In *International Conference on Learning Representations (ICLR)*, 2023b.
- Jinyoung Choi, Junoh Kang, and Bohyung Han. Enhanced diffusion sampling via extrapolation with multiple ode solutions. *arXiv preprint arXiv:2504.01855*, 2025.
- Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 248–255. IEEE, 2009. doi: 10.1109/CVPR.2009.5206848.

- 540 Carles Domingo-Enrich, Michal Drozdal, Brian Karrer, and Ricky TQ Chen. Adjoint matching:
541 Fine-tuning flow and diffusion generative models with memoryless stochastic optimal control.
542 *arXiv preprint arXiv:2409.08861*, 2024.
- 543
544 Lawrence C. Evans. *Partial Differential Equations*, volume 19 of *Graduate Studies in Mathematics*.
545 American Mathematical Society, Providence, RI, 2 edition, 2010.
- 546
547 Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,
548 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. DPOK: Reinforcement learning for
549 fine-tuning text-to-image diffusion models. In *Neurips*, volume 36, pp. 79858–79885, 2023.
- 550
551 Xuefeng Gao, Jiale Zha, and Xun Yu Zhou. Reward-directed score-based diffusion models via
552 q-learning. *arXiv preprint arXiv:2409.04832*, 2024.
- 553
554 Xuefeng Gao, Hoang M Nguyen, and Lingjiong Zhu. Wasserstein convergence guarantees for a
555 general class of score-based generative models. *J. Mach. Learn. Res.*, 26(43):1–54, 2025.
- 556
557 Yinbin Han, Meisam Razaviyayn, and Renyuan Xu. Stochastic control for fine-tuning diffusion
558 models: Optimality, regularity, and convergence. In *Proceedings of the 42nd International Con-*
559 *ference on Machine Learning*, 2025.
- 560
561 U. G. Haussmann and É. Pardoux. Time reversal of diffusions. *Ann. Probab.*, 14(4):1188–1205,
562 1986.
- 563
564 John R. Hershey and Peder A. Olsen. Approximating the kullback leibler divergence between gaussian
565 mixture models. In *2007 IEEE International Conference on Acoustics, Speech and Signal*
566 *Processing - ICASSP '07*, volume 4, pp. IV–317–IV–320, 2007. doi: 10.1109/ICASSP.2007.
567 366913.
- 568
569 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on*
570 *Deep Generative Models and Downstream Applications*, 2021.
- 571
572 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances*
573 *in Neural Information Processing Systems (NeurIPS)*, 2020.
- 574
575 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J.
576 Fleet. Video diffusion models. In *Advances in Neural Information Processing Systems (NeurIPS)*,
577 volume 35, 2022.
- 578
579 Daniel Zhengyu Huang, Jiaoyang Huang, and Zhengjiang Lin. Fast convergence for high-order ode
580 solvers in diffusion probabilistic models. *arXiv preprint arXiv:2506.13061*, 2025.
- 581
582 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-
583 a-pic: An open dataset of user preferences for text-to-image generation. In *Advances in Neural*
584 *Information Processing Systems (NeurIPS)*, 2023.
- 585
586 Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel,
587 Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human
588 feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- 589
590 Kyungmin Lee, Xiaohang Li, Qifei Wang, Junfeng He, Junjie Ke, Ming-Hsuan Yang, Irfan Essa,
591 Jinwoo Shin, Feng Yang, and Yinxiao Li. Calibrated multi-preference optimization for aligning
592 diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
593 *Recognition (CVPR)*, 2025.
- 594
595 Junzhe Li, Yutao Cui, Tao Huang, Yinping Ma, Chun Fan, Miles Yang, and Zhao Zhong. MixGRPO:
596 Unlocking flow-based GRPO efficiency with mixed ODE-SDE. *arXiv preprint arXiv:2507.21802*,
597 2025.
- 598
599 Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song, Mingda Wan, and Yufa Zhou. Unraveling
600 the smoothness properties of diffusion models: A gaussian mixture perspective. In *Proceedings*
601 *of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2025.

- 594 Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and
595 Mark D Plumbley. AudioLDM: Text-to-audio generation with latent diffusion models. In An-
596 dreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan
597 Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume
598 202 of *Proceedings of Machine Learning Research*, pp. 21450–21474. PMLR, 23–29 Jul 2023.
- 599 Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan,
600 Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *arXiv*
601 *preprint arXiv:2505.05470*, 2025.
- 602 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*
603 *ence on Learning Representations (ICLR)*, 2019.
- 604 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast
605 ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in neural*
606 *information processing systems*, 35:5775–5787, 2022.
- 607 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast
608 solver for guided sampling of diffusion probabilistic models. *Machine Intelligence Research*, pp.
609 1–22, 2025.
- 610 Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte carlo gradient esti-
611 mation in machine learning. *Journal of Machine Learning Research*, 21(132):1–62, 2020.
- 612 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
613 Zhang, Sandhini Agarwal, Katarina Slama, and Alex Ray. Training language models to follow
614 instructions with human feedback. In *Neurips*, volume 35, pp. 27730–27744, 2022.
- 615 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
616 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
617 synthesis. In *International Conference on Learning Representations (ICLR)*, 2024.
- 618 Hannes Risken. *Fokker-Planck Equation*, pp. 63–95. Springer Berlin Heidelberg, Berlin, Heidel-
619 berg, 1996. ISBN 978-3-642-61544-3. doi: 10.1007/978-3-642-61544-3_4.
- 620 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
621 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Con-*
622 *ference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- 623 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
624 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,
625 Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev.
626 Laion-5b: An open large-scale dataset for training next generation image-text models. In *Neurips*,
627 volume 35, pp. 25278–25294, 2022.
- 628 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
629 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 630 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
631 Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathe-
632 matical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 633 Joar Skalse, Nikolaus H. R. Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and
634 characterizing reward hacking. *arXiv preprint arXiv:2209.13085*, 2022.
- 635 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*,
636 2021a.
- 637 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
638 Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*,
639 2021b.
- 640 Dmitrii Sorokin, Maksim Nakhodnov, Andrey Kuznetsov, and Aibek Alanov. Imagereff: Balancing
641 quality and diversity in human-aligned diffusion models. *arXiv preprint arXiv:2505.22569*, 2025.

- 648 Daniel W. Stroock and S. R. Srinivasa Varadhan. *Multidimensional diffusion processes*, volume 233
649 of *Grundlehren der Mathematischen Wissenschaften*. Springer-Verlag, 1979.
- 650
- 651 Wenpin Tang. Fine-tuning of diffusion models via stochastic control: entropy regularization and
652 beyond. *arXiv preprint arXiv:2403.06279*, 2024.
- 653
- 654 Wenpin Tang and Hanyang Zhao. Contractive diffusion probabilistic models. *arXiv preprint*
655 *arXiv:2401.13115*, 2024.
- 656
- 657 Wenpin Tang and Hanyang Zhao. Score-based diffusion models via stochastic differential equations.
658 *Stat. Surv.*, 19:28–64, 2025.
- 659
- 660 Masatoshi Uehara, Yulai Zhao, Kevin Black, Ehsan Hajiramezani, Gabriele Scalia, Nathaniel Lee
661 Diamant, Alex M Tseng, Tommaso Biancalani, and Sergey Levine. Fine-tuning of continuous-
662 time diffusion models as entropy-regularized control. *arXiv preprint arXiv:2402.15194*, 2024.
- 663
- 664 Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Comput.*,
665 23(7):1661–1674, 2011.
- 666
- 667 Yibin Wang, Yuhang Zang, Hao Li, Cheng Jin, and Jiaqi Wang. Unified reward model for multi-
668 modal understanding and generation. *arXiv preprint arXiv:2503.05236*, 2025.
- 669
- 670 Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
671 learning. *Mach. Learn.*, 8(3-4):229–256, 1992.
- 672
- 673 Genta Indra Winata, Hanyang Zhao, Anirban Das, Wenpin Tang, David D Yao, Shi-Xiong Zhang,
674 and Sambit Sahu. Preference tuning with human feedback on language, speech, and vision tasks:
675 A survey. *J. Artif. Intell. Res.*, 82:2595–2661, 2025.
- 676
- 677 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.
678 Human preference score v2: A solid benchmark for evaluating human preferences of text-to-
679 image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.
- 680
- 681 Yuchen Wu, Minshuo Chen, Zihao Li, Mengdi Wang, and Yuting Wei. Theoretical insights for
682 diffusion guidance: A case study for Gaussian mixture models. In Ruslan Salakhutdinov, Zico
683 Kolter, Katharine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp
684 (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of
685 *Proceedings of Machine Learning Research*, pp. 53291–53327. PMLR, 21–27 Jul 2024.
- 686
- 687 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
688 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.
689 In *Neurips*, 2023.
- 690
- 691 Zeyue Xue, Jie Wu, Yu Gao, Fangyuan Kong, Lingting Zhu, Mengzhao Chen, Zhiheng Liu, Wei Liu,
692 Qiushan Guo, Weilin Huang, and Ping Luo. Dancegrp: Unleashing GRPO on visual generation.
693 *arXiv preprint arXiv:2505.07818*, 2025.
- 694
- 695 Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiabin Chen, Weihang Shen, Xiaolong Zhu, and Xiu Li.
696 Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings*
697 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8941–8951, 2024.
- 698
- 699 Oussama Zekri and Nicolas Boullé. Fine-tuning discrete diffusion models with policy gradient
700 methods. *arXiv preprint arXiv:2502.01384*, 2025.
- 701
- 702 Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential integrator.
703 *arXiv preprint arXiv:2204.13902*, 2022.
- 704
- 705 Qinsheng Zhang, Molei Tao, and Yongxin Chen. gDDIM: generalized denoising diffusion implicit
706 models. In *ICLR*, 2023.
- 707
- 708 Yinan Zhang, Eric Tzeng, Yilun Du, and Dmitry Kislyuk. Large-scale reinforcement learning for
709 diffusion models. In *European Conference on Computer Vision*, pp. 1–17. Springer, 2024a.

702 Ziyi Zhang, Sen Zhang, Yibing Zhan, Yong Luo, Yonggang Wen, and Dacheng Tao. Confronting
703 reward over-optimization for diffusion models: A perspective of inductive and primacy biases.
704 *arXiv preprint arXiv:2402.08552*, 2024b.
705
706 Hanyang Zhao, Haoxian Chen, Ji Zhang, David Yao, and Wenpin Tang. Scores as Actions: a
707 framework of fine-tuning diffusion models by continuous-time reinforcement learning. *arXiv*
708 *preprint arXiv:2409.08400*, 2024.
709 Hanyang Zhao, Haoxian Chen, Ji Zhang, David Yao, and Wenpin Tang. Score as Action: Fine
710 tuning diffusion generative models by continuous-time reinforcement learning. In *ICML*, 2025.
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A LEMMAS ON LINEAR DYNAMICS WITH GAUSSIAN PRIORS

Lemma A.1. A stochastic process $\{Z_t\}_{t=0}^T$ with first order linear dynamic and initial Gaussian distribution

$$\begin{cases} dZ_t = f(t)Z_t dt + g(t)dt + h(t)dB_t & t \in [0, T] \\ Z_0 \sim \mathcal{N}(0, 1) \end{cases}$$

is distributed following

$$Z_t \sim \mathcal{N}\left(e^{F(t)} \int_0^t e^{-F(s)} g(s) ds, e^{2F(t)} \int_0^t e^{-2F(s)} h^2(s) ds\right),$$

in which $F(t)$ is the integrating factor satisfying $F(t) = \int_0^t f(s) ds$.

Lemma A.2. A stochastic process $\{Z_t\}_{t=0}^T$ with first order linear dynamic and initial Gaussian distribution

$$\begin{cases} dZ_t = f(t)Z_t dt + g(t)dB_t & t \in [0, T] \\ Z_0 \sim \mathcal{N}(\mu_Z, \sigma_Z^2) \end{cases}$$

is distributed following

$$\begin{aligned} Z_t &\sim Z_0 \cdot e^{F(t)} + \mathcal{N}\left(0, \int_0^t e^{-2F(s)} g^2(s) ds\right) \cdot e^{F(t)} \\ &\sim \mathcal{N}\left(\mu_Z \cdot e^{F(t)}, \left[\sigma_Z^2 + \int_0^t e^{-2F(s)} g^2(s) ds\right] \cdot e^{2F(t)}\right), \end{aligned}$$

in which $F(t)$ is the integrating factor satisfying $F(t) = \int_0^t f(s) ds$.

Lemma A.3. A parametrized family of stochastic processes $\{Z_t^\theta\}_{t=0}^T$ with initial Gaussian distribution

$$\begin{cases} dZ_t^\theta = f(t) \cdot (Z_t^\theta + \theta(t)) dt + g(t) dB_t & t \in [0, T] \\ Z_0 \sim \mathcal{N}(\mu_Z, \sigma_Z^2) \end{cases}$$

is distributed following

$$\begin{aligned} Z_t &\sim Z_0 \cdot e^{F(t)} + \mathcal{N}\left(\int_0^t e^{-F(s)} f(s) \theta(s) ds, \int_0^t e^{-2F(s)} g^2(s) ds\right) \cdot e^{F(t)} \\ &\sim \mathcal{N}\left(\mu_Z \cdot e^{F(t)} + \int_0^t e^{-F(s)} f(s) \theta(s) ds, \left[\sigma_Z^2 + \int_0^t e^{-2F(s)} g^2(s) ds\right] \cdot e^{2F(t)}\right), \end{aligned}$$

in which $F(t)$ is the integrating factor satisfying $F(t) = \int_0^t f(s) ds$.

B USEFUL PROPOSITIONS

Proposition B.1. *The terminal distribution of Variance Exploding parametrized backward dynamics $\{Y_t^\theta\}_{t=0}^T$ is always Gaussian following the law*

$$\mathcal{N}(\mu_{Y_T^\theta}, \sigma_{Y_T^\theta}^2) := \mathcal{N}\left(\theta \cdot (1 + T^2)^{-\frac{1+\eta^2}{2}} - \theta, 1 - (1 + T^2)^{-(1+\eta^2)}\right).$$

Therefore, reward function (11) takes the form of

$$\mathcal{J}_\eta(\theta) = -\left(\sigma_{Y_T^\theta}^2 + (1 - \mu_{Y_T^\theta})^2\right).$$

Proof. By comparing the coefficients with Lemma A.1, we have

$$\begin{cases} f_\eta(t) = -\frac{(1+\eta^2)(T-t)}{1+(T-t)^2} \\ g_\eta(t) = \eta\sqrt{2(T-t)} \\ \sigma_Z^2 = T^2. \end{cases}$$

Therefore, we first examine the exponential of integrating factor,

$$\begin{aligned} e^{F_\eta(t)} &= \exp\left(\int_0^t f_\eta(s) ds\right) \\ &= \exp\left(\int_0^t -\frac{(1+\eta^2)(T-s)}{1+(T-s)^2} ds\right) \\ &= \exp\left(\int_{1+T^2}^{1+(T-t)^2} \frac{(1+\eta^2)d(1+(T-s)^2)}{2(1+(T-s)^2)}\right) \\ &= \left(\frac{1+T^2}{1+(T-t)^2}\right)^{-\frac{1+\eta^2}{2}}. \end{aligned}$$

At terminal time T , the cumulative factor is

$$e^{F_\eta(T)} = \exp\left(\int_0^T f_\eta(s) ds\right) = (1+T^2)^{-\frac{1+\eta^2}{2}}.$$

Also, the cumulative Gaussian variance generated from the backward process is,

$$\begin{aligned} \int_0^t e^{-2F_\eta(s)} g_\eta^2(s) ds &= \int_0^t \left(\frac{1+T^2}{1+(T-s)^2}\right)^{(1+\eta^2)} \cdot (2\eta^2(T-s)) ds \\ &= \int_{1+T^2}^{1+(T-t)^2} (-\eta^2) \left(\frac{1+T^2}{1+(T-s)^2}\right)^{(1+\eta^2)} d(1+(T-s)^2) \\ &= (1+T^2)^{(1+\eta^2)} \int_{1+T^2}^{1+(T-t)^2} (-\eta^2)(1+(T-s)^2)^{-(1+\eta^2)} d(1+(T-s)^2) \\ &= (1+T^2)^{(1+\eta^2)} \cdot \left((1+(T-t)^2)^{-\eta^2} - (1+T^2)^{-\eta^2}\right). \end{aligned}$$

At terminal time T , the variance from the process is

$$\begin{aligned} \int_0^T e^{-2F_\eta(s)} g_\eta^2(s) ds &= (1+T^2)^{(1+\eta^2)} \cdot \left(1 - (1+T^2)^{-\eta^2}\right) \\ &= (1+T^2)^{(1+\eta^2)} - (1+T^2). \end{aligned}$$

Together with the initial Gaussian variance, the terminal distribution remains a zero-mean Gaussian:

$$\begin{aligned} Y_T &\sim \mathcal{N}\left(0, [(1+T^2)^{(1+\eta^2)} - (1+T^2 - \sigma_Z)] \cdot (1+T^2)^{-(1+\eta^2)}\right) \\ &\sim \mathcal{N}\left(0, [(1+T^2)^{(1+\eta^2)} - 1] \cdot (1+T^2)^{-(1+\eta^2)}\right) \\ &\sim \mathcal{N}\left(0, 1 - (1+T^2)^{-(1+\eta^2)}\right). \end{aligned}$$

Now we consider the terminal distribution for the parametrized process Y_T^θ . To reduce the problem to a dynamic with linear drift term, we define

$$Z_t^\theta = Y_t^\theta + \theta.$$

so that

$$\begin{cases} dZ_t = f_\eta(t)Z_t dt + g_\eta(t)dB_t & t \in [0, T] \\ Z_0 \sim \mathcal{N}(\theta, T^2) \end{cases}$$

Observe that both the dynamic variance and the initial distribution variance for $\{Z_t^\theta\}$ and $\{Y_t\}$ are the same, we may directly apply Lemma A.1 to obtain

$$Z_T^\theta \sim \mathcal{N}(\theta \cdot (1 + T^2)^{-\frac{1+\eta^2}{2}}, 1 - (1 + T^2)^{-(1+\eta^2)})$$

Therefore,

$$Y_T^\theta \sim \mathcal{N}(\theta \cdot (1 + T^2)^{-\frac{1+\eta^2}{2}} - \theta, 1 - (1 + T^2)^{-(1+\eta^2)}),$$

and thus

$$\begin{cases} \mu_{Y_T^\theta} := \theta \cdot (1 + T^2)^{-\frac{1+\eta^2}{2}} - \theta \\ \sigma_{Y_T^\theta}^2 := 1 - (1 + T^2)^{-(1+\eta^2)} \end{cases} \quad (16)$$

In addition, we can now give a closed form representation of our reward,

$$\begin{aligned} \mathcal{J}_\eta(\theta) &= \mathbb{E}[(Y_T^\theta - 1)^2] \\ &= \mathbb{E}[Y_T^2] - 2\mathbb{E}[Y_T] + 1 \\ &= \sigma_{Y_T^\theta}^2 + \mu_{Y_T^\theta}^2 - 2\mu_{Y_T^\theta} + 1 \\ &= \sigma_{Y_T^\theta}^2 + (1 - \mu_{Y_T^\theta})^2, \end{aligned}$$

which yields to the desired expression. \square

Proposition B.2. Given β, η, T , the unique maximizer to the entropy regularized target (8) is:

$$\theta_\eta^* = -\left(1 + \frac{\beta}{2}\right) \cdot \left[1 - (1 + T^2)^{-\frac{1+\eta^2}{2}}\right]^{-1}.$$

Proof. A classical distance result on two Gaussian distributions (Hershey & Olsen, 2007) states:

Lemma B.1. The KL divergence for two Gaussian distributions $P \sim \mathcal{N}(\mu_1, \sigma_1)$ and $Q \sim \mathcal{N}(\mu_2, \sigma_2)$,

$$\text{KL}(P||Q) = \log\left(\frac{\sigma_2}{\sigma_1}\right) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}.$$

In our model, $P \sim Y_T^\theta$ and $Q \sim Y_T^\odot$, so $\mu_1 = \mu_{Y_T^\theta}$, $\sigma_1 = \sigma_{Y_T}$, $\mu_2 = 0$, $\sigma_2 = 1$.

$$\begin{aligned} \text{KL}(Y_T^\theta||Y_T^\odot) &= \log\left(\frac{1}{\sigma_{Y_T}}\right) + \frac{\sigma_{Y_T}^2 + \mu_{Y_T^\theta}^2 - 1}{2} \\ &= -\log(\sigma_{Y_T}) + \frac{\sigma_{Y_T}^2 + \mu_{Y_T^\theta}^2 - 1}{2}. \end{aligned}$$

Therefore,

$$\begin{aligned} -F_\eta(\theta) &= \sigma_{Y_T}^2 + (\mu_{Y_T^\theta} - 1)^2 + \left(-\beta \log(\sigma_{Y_T}) + \frac{\beta}{2} \cdot (\sigma_{Y_T}^2 + \mu_{Y_T^\theta}^2 - 1)\right) \\ &= (\sigma_{Y_T}^2 - \beta \log(\sigma_{Y_T}) + \frac{\beta}{2} \sigma_{Y_T}^2 - \frac{\beta}{2}) + (\mu_{Y_T^\theta} - 1)^2 + \frac{\beta}{2} \mu_{Y_T^\theta}^2 \\ &= (\sigma_{Y_T}^2 - \beta \log(\sigma_{Y_T}) + \frac{\beta}{2} \sigma_{Y_T}^2 - \frac{\beta}{2} + 1) + \left(1 + \frac{\beta}{2}\right) \mu_{Y_T^\theta}^2 - 2\mu_{Y_T^\theta}. \end{aligned}$$

So it suffices to minimize a quadratic function w.r.t $\mu_{Y_T^\theta}$, of which we know

$$\mu_{Y_T^{\theta^*}} = -\frac{-2}{2(1 + \frac{\beta}{2})} = (1 + \frac{\beta}{2})^{-1};$$

and thus

$$\theta_\eta^* \cdot [(1 + T^2)^{-\frac{1+\eta^2}{2}} - 1] = \mu_{Y_T^{\theta^*}} = (1 + \frac{\beta}{2})^{-1}.$$

This gives us the unique maximizer

$$\theta_\eta^* = \left((1 + \frac{\beta}{2}) \cdot [(1 + T^2)^{-\frac{1+\eta^2}{2}} - 1] \right)^{-1}$$

as desired. \square

Proposition B.3. *The terminal distribution of Variance Preserving parametrized backward dynamics $\{Y_t^\theta\}_{t=0}^T$ is always Gaussian following the law*

$$\mathcal{N}(\mu_{Y_T^\theta}, 1) := \mathcal{N}\left(\theta \cdot e^{-\frac{(1+\eta^2) \cdot T^2}{2}} - \theta, 1\right).$$

And reward function (11) takes the form of $\mathcal{J}_\eta(\theta) = -(1 + (1 - \mu_{Y_T^\theta})^2)$.

Proof. By comparing with the coefficients in Lemma A.2, we have

$$\begin{cases} f_\eta(t) = -\eta^2(T-t) \\ g_\eta(t) = -(1 + \eta^2)(T-t)e^{-\frac{(T-t)^2}{2}}\theta_\eta^* \\ h_\eta(t) = \eta\sqrt{2(T-t)} \end{cases}$$

Therefore, we first examine the exponential of integrating factor,

$$\begin{aligned} e^{F_\eta(t)} &= \exp\left(\int_0^t f_\eta(s)ds\right) \\ &= \exp\left(-\eta^2 \int_{T-t}^T s ds\right) \\ &= \exp\left(-\frac{\eta^2}{2} \cdot (T^2 - (T-t)^2)\right). \end{aligned}$$

At terminal time T , the cumulative factor is

$$e^{F_\eta(T)} = \exp\left(-\frac{\eta^2}{2} \cdot (T^2 - (T-T)^2)\right) = e^{-\frac{\eta^2 T^2}{2}}.$$

Now we are able to compute

$$\begin{aligned} \int_0^t e^{-F_\eta(s)} g(s) ds &= \int_0^t e^{\frac{\eta^2}{2}(T^2 - (T-s)^2)} \cdot \left(- (1 + \eta^2)(T-s)e^{-\frac{(T-s)^2}{2}}\theta_\eta^*\right) ds \\ &= (1 + \eta^2) \cdot \theta_\eta^* \cdot \int_{T-t}^T e^{\frac{\eta^2}{2}(T^2 - s^2)} \cdot \left(se^{-\frac{s^2}{2}}\right) ds \\ &= (1 + \eta^2) \cdot \theta_\eta^* \cdot \int_{T-t}^T e^{\frac{\eta^2 T^2}{2}} \cdot \left(se^{-\frac{(1+\eta^2)s^2}{2}}\right) ds \\ &= (1 + \eta^2) \cdot \theta_\eta^* \cdot e^{\frac{\eta^2 T^2}{2}} \cdot \left[\frac{1}{1 + \eta^2} \cdot e^{-\frac{(1+\eta^2)s^2}{2}}\right]_{s=T-t}^{s=T} \\ &= \theta_\eta^* \cdot e^{\frac{\eta^2 T^2}{2}} \cdot \left(e^{-\frac{(1+\eta^2)T^2}{2}} - e^{-\frac{(1+\eta^2)(T-t)^2}{2}}\right) \end{aligned}$$

972 Therefore,

$$973 \mu_{Y_T^\theta} = e^{-\frac{\eta^2 T^2}{2}} \cdot \theta_\eta^* \cdot e^{\frac{\eta^2 T^2}{2}} \cdot (e^{-\frac{(1+\eta^2)T^2}{2}} - e^0) = \theta_\eta^* \cdot (e^{-\frac{(1+\eta^2)T^2}{2}} - 1)$$

974 To the variance preserving property, it suffices to show

$$975 \frac{1}{e^{2F_\eta(t)}} = \int_0^t e^{-2F_\eta(s)} h_\eta^2(s) ds.$$

976 In fact,

$$977 \frac{d}{dt} e^{-2F_\eta(t)} = e^{-2F_\eta(t)} \cdot \frac{d(-\eta^2(T-t)^2)}{dt} = e^{-2F_\eta(t)} \cdot (\eta^2 \cdot 2(T-t)) = e^{-2F_\eta(t)} \cdot h_\eta^2(t).$$

978 □

979 **Proposition B.4.** Given β, η, T , the unique maximizer to the entropy regularized target (8) is:

$$980 \theta_\eta^* = -\left(\left(1 + \frac{\beta}{2}\right) \cdot \left[1 - e^{-\frac{(1+\eta^2) \cdot T^2}{2}}\right] \right)^{-1}.$$

981 *Proof.* Since $\sigma \equiv 1$, according to Lemma B.1, the maximum reward is attained at

$$982 \mu_{Y_T^{\theta^*}} = -\frac{-2}{2\left(1 + \frac{\beta}{2}\right)} = \left(1 + \frac{\beta}{2}\right)^{-1}.$$

983 By Proposition B.1,

$$984 \theta_\eta^* \cdot \left[e^{-\frac{(1+\eta^2) \cdot T^2}{2}} - 1 \right] = \mu_{Y_T^{\theta^*}} = \left(1 + \frac{\beta}{2}\right)^{-1}.$$

985 This gives us the unique maximizer

$$986 \theta_\eta^* = \left(\left(1 + \frac{\beta}{2}\right) \cdot \left[e^{-\frac{(1+\eta^2) \cdot T^2}{2}} - 1 \right] \right)^{-1}.$$

987 □

988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

C PROOFS FOR SECTION 3

C.1 PROOF OF THEOREM 3.1

We first consider the Y_T^{ODE} process as discussed in Section 3.1. Since $\eta = 0$, by Proposition B.2,

$$\mu_{Y_T^{ODE}} = \theta_\eta^* \cdot (1 + T^2)^{-\frac{1}{2}} - \theta_\eta^* = (1 + \frac{\beta}{2})^{-1} \cdot \frac{(1 + T^2)^{-\frac{1}{2}} - 1}{(1 + T^2)^{-\frac{1+\eta^2}{2}} - 1},$$

and

$$\sigma_{Y_T^{ODE}}^2 = \sigma_{Y_T^0}^2 = 1 - (1 + T^2)^{-1}.$$

Moreover, the quadratic reward \mathcal{J}_{ODE} for Y_T^{ODE} is

$$\begin{aligned} \mathcal{J}_0(\theta_\eta^*) &= -(\sigma_{Y_T^{ODE}}^2 + (\mu_{Y_T^{ODE}} - 1)^2) \\ &= (-1) + (1 + T^2)^{-1} - \left(1 - (1 + \frac{\beta}{2})^{-1} \cdot \frac{(1 + T^2)^{-\frac{1}{2}} - 1}{(1 + T^2)^{-\frac{1+\eta^2}{2}} - 1}\right)^2. \end{aligned}$$

Similarly, reward \mathcal{J}_{SDE} for Y_T^{SDE} is

$$\begin{aligned} \mathcal{J}_\eta(\theta_\eta^*) &= -(\sigma_{Y_T^{SDE}}^2 + (\mu_{Y_T^{SDE}} - 1)^2) \\ &= (-1) + (1 + T^2)^{-(1+\eta^2)} - \left(1 - (1 + \frac{\beta}{2})^{-1}\right)^2. \end{aligned}$$

For simplification, we denote

$$\bar{T} := 1 + T^2 \in (T^2, 2T^2), \quad \bar{\beta} := (1 + \frac{\beta}{2})^{-1} \in (0, 1].$$

Now the reward gap

$$\begin{aligned} \Delta_\eta &= \mathcal{J}_{SDE} - \mathcal{J}_{ODE} \\ &= \left((1 + T^2)^{-(1+\eta^2)} - (1 + T^2)^{-1} \right) \\ &\quad - \left(\left(1 - (1 + \frac{\beta}{2})^{-1}\right)^2 - \left(1 - (1 + \frac{\beta}{2})^{-1} \cdot \frac{(1 + T^2)^{-\frac{1}{2}} - 1}{(1 + T^2)^{-\frac{1+\eta^2}{2}} - 1}\right)^2 \right) \\ &= \left(\frac{\bar{T}^{-\eta^2} - 1}{\bar{T}} \right) - \left((1 - \bar{\beta})^2 - \left(1 - \bar{\beta} \cdot \frac{\bar{T}^{-\frac{1}{2}} - 1}{\bar{T}^{-\frac{1+\eta^2}{2}} - 1}\right)^2 \right). \end{aligned}$$

Since $0 < \bar{T}^{-\eta^2} \leq 1$, we can bound

$$-\bar{T}^{-1} < \frac{\bar{T}^{-\eta^2} - 1}{\bar{T}} \leq 0,$$

and

$$1 - \bar{T}^{-\frac{1}{2}} = \frac{\bar{T}^{-\frac{1}{2}} - 1}{0 - 1} < \frac{\bar{T}^{-\frac{1}{2}} - 1}{\bar{T}^{-\frac{1+\eta^2}{2}} - 1} \leq \frac{\bar{T}^{-\frac{1}{2}} - 1}{\bar{T}^{-\frac{1}{2}} - 1} = 1.$$

Therefore,

$$\begin{aligned} |\Delta_\eta| &= \left| \frac{\bar{T}^{-\eta^2} - 1}{\bar{T}} \right| + \left| (1 - \bar{\beta})^2 - \left(1 - \bar{\beta} \cdot \frac{\bar{T}^{-\frac{1}{2}} - 1}{\bar{T}^{-\frac{1+\eta^2}{2}} - 1}\right)^2 \right| \\ &\leq \bar{T}^{-1} + \left| (1 - \bar{\beta})^2 - \left((1 - \bar{\beta}) + \bar{\beta} \cdot \bar{T}^{-\frac{1}{2}} \right)^2 \right| \\ &\leq \bar{T}^{-1} + \left(\bar{\beta}^2 \cdot \bar{T}^{-1} + 2\bar{\beta}(1 - \bar{\beta})\bar{T}^{-\frac{1}{2}} \right) \\ &\leq 2\bar{\beta}(1 - \bar{\beta})\bar{T}^{-1} + o(\bar{T}^{-1}) \\ &\leq \frac{1}{2\bar{T}} + o\left(\frac{1}{\bar{T}}\right). \end{aligned}$$

1080 Now we attempt to bound the reference gap. The reward \mathcal{J}_{REF} for Y_T^{REF} is

$$1081 \mathcal{J}_\eta(0) = -(\sigma_{Y_T^{REF}}^2 + (\mu_{Y_T^{REF}} - 1)^2)$$

$$1082 = (-1) + (1 + T^2)^{-(1+\eta^2)} - 1.$$

1083 Finally we have,

$$1084 \Delta_\eta^{REF} = \mathcal{J}_{ODE} - \mathcal{J}_{REF}$$

$$1085 = \bar{T}^{-1} - \left(1 - \bar{\beta} \cdot \frac{\bar{T}^{-\frac{1}{2}} - 1}{\bar{T}^{-\frac{1+\eta^2}{2}} - 1}\right)^2 - (\bar{T}^{-(1+\eta^2)} - 1)$$

$$1086 \geq \bar{T}^{-1} - \left(2\bar{\beta}(1 - \bar{\beta})\bar{T}^{-\frac{1}{2}} + o(\bar{T}^{-\frac{1}{2}})\right) - (\bar{T}^{-1} - 1)$$

$$1087 \geq 1 - \frac{1}{2\bar{T}} + o\left(\frac{1}{\bar{T}}\right)$$

1097 C.2 PROOF OF THEOREM 3.2

1098 Similar to Appendix C.1, the quadratic reward (11) for Y_T^{ODE} is

$$1099 \mathcal{J}_{ODE} = -(\sigma_{Y_T^{ODE}}^2 + (\mu_{Y_T^{ODE}} - 1)^2)$$

$$1100 = (-1) - \left(1 - \left(1 + \frac{\beta}{2}\right)^{-1} \cdot \frac{e^{-\frac{T^2}{2}} - 1}{e^{-\frac{(1+\eta^2)T^2}{2}} - 1}\right)^2.$$

1101 And the reward for Y_T^{SDE} is

$$1102 \mathcal{J}_{SDE} = -(\sigma_{Y_T^{SDE}}^2 + (\mu_{Y_T^{SDE}} - 1)^2)$$

$$1103 = (-1) - \left(1 - \left(1 + \frac{\beta}{2}\right)^{-1}\right)^2.$$

1104 With $\bar{\beta} := (1 + \frac{\beta}{2})^{-1} \in (0, 1]$,

$$1105 |\Delta_\eta| = \left| (1 - \bar{\beta})^2 - \left(1 - \bar{\beta} \cdot \frac{e^{-\frac{T^2}{2}} - 1}{e^{-\frac{(1+\eta^2)T^2}{2}} - 1}\right)^2 \right|$$

$$1106 \leq \left| (1 - \bar{\beta})^2 - \left((1 - \bar{\beta}) + \bar{\beta} \cdot e^{-\frac{T^2}{2}}\right)^2 \right|$$

$$1107 \leq \frac{e^{-\frac{T^2}{2}}}{2} + o(e^{-T^2}).$$

1108 Finally, $Y_T^{REF} \sim \mathcal{N}(0, 1)$, so

$$1109 \mathcal{J}_{REF} = (-1) - (1 - 0)^2.$$

1110 And thus

$$1111 \Delta_\eta^{REF} = \mathcal{J}_{ODE} - \mathcal{J}_{REF}$$

$$1112 = 1 - \left(1 - \left(1 + \frac{\beta}{2}\right)^{-1} \cdot \frac{e^{-\frac{T^2}{2}} - 1}{e^{-\frac{(1+\eta^2)T^2}{2}} - 1}\right)^2$$

$$1113 \geq 1 - \frac{e^{-\frac{T^2}{2}}}{2} + o(e^{-T^2})$$

C.3 PROOF OF COROLLARY 3.1

Let $\bullet \in \{SDE, ODE\}$. We decompose $Y^{\theta, \bullet} := Y_{\parallel}^{\theta, \bullet} + Y_{\perp}^{\theta, \bullet}$ in terms of \mathbf{r} . Therefore,

$$r(Y^{\theta, \bullet}(T)) = -\left|Y_{\perp}^{\theta, \bullet} + (Y_{\parallel}^{\theta, \bullet} - \mathbf{r})\right|^2 = -\left(\left|Y_{\perp}^{\theta, \bullet}\right|^2 + \left|Y_{\parallel}^{\theta, \bullet} - \mathbf{r}\right|^2\right)$$

Since $\mu_i \perp \mathbf{r}$, $\mathbb{E}[Y_{\parallel}^{\theta}(0)] = 0$. Also, the backward dynamic injects a scaled multiple of \mathbf{I}_d noise, so the coordinate-wise dynamics are independent. Therefore, we are able to analyze the $\text{Span}(\mathbf{r})$ subspace via separating each dimension $e_k \in \text{Span}(\mathbf{r})$, in which $\left\{\text{Proj}_{e_k}(Y_{\parallel}^{\theta, \bullet, SDE})\right\}$ and $\left\{\text{Proj}_{e_k}(Y_{\parallel}^{\theta, \bullet, ODE})\right\}$ follows a similar controlled motion as discussed in Theorem 3.1 and 3.2. Moreover, the score function can be bounded by a constant of $\max_i \sigma_i, \min_i^{-1} \sigma_i$ (Liang et al., 2025).

On the other hand, $\theta_{\perp} = \mathbf{0}$ is a minimizer for $\left|Y_{\perp}^{\theta, \bullet}\right|^2$, since an additional drift perpendicular to \mathbf{r} does not alter reward variance but pushes away reward mean of $Y_{\perp}^{\theta, \bullet}(T)$ from reference priors. Therefore, $\theta = \theta_{\parallel}$. A similar analysis on Gaussian priors with controlled drifts for VP and VE dynamics yields to the desired bounds.

D PROOF FOR SECTION 4

Proposition D.1. (Gronwall's) Suppose integrable functions $u, \alpha, \beta : [0, T] \rightarrow \mathbb{R}$ satisfies $u'(t) \leq \alpha(t)u(t) + \beta(t)$,

$$u(T) \leq e^{\int_0^T \alpha(s)ds} \left(u(0) + \int_0^T e^{-\int_0^s \alpha(\tau)d\tau} \beta(s)ds \right). \quad (17)$$

A proof can be found in (Evans, 2010).

D.1 PROOF OF THEOREM 4.1

Let $u(t) := \mathbb{E}[\|Y_t^{ODE} - Y_t^{SDE}\|^2]$. Note that $u(0) = 0$. Therefore, with appropriate α, β satisfying the conditions of equation 17,

$$u(T) \leq e^{\alpha T} \left(\int_0^T \beta \cdot e^{-\alpha s} ds \right) = \frac{(e^{\alpha T} - 1) \cdot \beta}{\alpha}. \quad (18)$$

By Ito's Lemma,

$$\begin{aligned} u'(t) &= \frac{d}{dt} \mathbb{E}[\|Y_t^{ODE} - Y_t^{SDE}\|^2] \\ &= 2\mathbb{E}\langle Y_t^{ODE} - Y_t^{SDE}, -f(t, Y_t^{ODE}) + f(t, Y_t^{SDE}) \rangle \\ &\quad + g^2(t)\mathbb{E}\langle Y_t^{ODE} - Y_t^{SDE}, s_{\theta}(t, Y_t^{ODE}) - s_{\theta}(t, Y_t^{SDE}) \rangle \\ &\quad - g^2(t)\mathbb{E}\langle Y_t^{ODE} - Y_t^{SDE}, \eta^2 \cdot s_{\theta}(t, Y_t^{SDE}) \rangle + \eta^2 g^2(t). \end{aligned}$$

According to (Tang & Zhao, 2025),

$$\mathbb{E}\langle Y_t^{ODE} - Y_t^{SDE}, -f(t, Y_t^{ODE}) + f(t, Y_t^{SDE}) \rangle = \begin{cases} 0 & \text{for VE,} \\ g^2(t)u(t)/2 & \text{for VP.} \end{cases}$$

By Condition 1,

$$\begin{aligned} \mathbb{E}\langle Y_t^{ODE} - Y_t^{SDE}, s_{\theta}(t, Y_t^{ODE}) - s_{\theta}(t, Y_t^{SDE}) \rangle &\leq \mathbb{E}\langle Y_t^{ODE} - Y_t^{SDE}, L \cdot (Y_t^{ODE} - Y_t^{SDE}) \rangle \\ &= L \cdot u(t). \end{aligned}$$

By Condition 2,

$$\begin{aligned} |\mathbb{E}\langle Y_t^{\text{ODE}} - Y_t^{\text{SDE}}, \eta^2 \cdot s_\theta(t, Y_t^{\text{SDE}}) \rangle| &\leq u(t) + \frac{\eta^4}{4} \cdot \mathbb{E}[\|s_\theta(t, Y_t^{\text{SDE}})\|^2] \\ &= u(t) + \frac{\eta^4}{4} \cdot A. \end{aligned}$$

Together, by absorbing $\|g\|_\infty$ into an η -free constant C_{scheme} ,

$$\begin{aligned} u'(t) &\leq \|g\|_\infty^2 u(t) + \|g\|_\infty^2 \cdot L \cdot u(t) + \|g\|_\infty^2 \left(u(t) + \frac{\eta^4}{4} A \right) + \|g\|_\infty^2 \cdot \eta^2 \\ &\leq \|g\|_\infty^2 u(t) + \|g\|_\infty^2 \cdot L \cdot u(t) + \|g\|_\infty^2 \left(u(t) + \frac{\eta^4}{4} A \right) + \|g\|_\infty^2 \cdot \eta^2 \\ &\leq \underbrace{(C_{\text{scheme}} \cdot L)}_\alpha \cdot u(t) + \underbrace{\max\{\eta^4, \eta^2\} \cdot \left(\frac{A}{4} + 1 \right)}_\beta \cdot \|g\|_\infty^2. \end{aligned}$$

And we apply equation 18 with $T = 1$:

$$\begin{aligned} u(T) &\leq \frac{(e^{\alpha T} - 1) \cdot \beta}{\alpha} \\ &\leq \frac{\exp(C_{\text{scheme}} \cdot L)}{C_{\text{scheme}} \cdot L} \cdot \max\{\eta^4, \eta^2\} \cdot \left(\frac{A}{4} + 1 \right). \end{aligned}$$

Finally, the W_2 distance can be bounded in terms of this L^2 -distance via the chosen coupling:

$$W_2(\mathcal{L}(Y_T^{\text{ODE}}), \mathcal{L}(Y_T^{\text{SDE}})) \leq \left(\mathbb{E}[\|Y_T^{\text{ODE}} - Y_T^{\text{SDE}}\|^2] \right)^{1/2} \leq C(\|g\|_\infty, L) \max\{\eta, \eta^2\} \sqrt{1 + A}.$$

D.2 PROOF FOR STRONGLY LOG-CONCAVE DISTRIBUTIONS

Assumption D.1. (Strong log-concavity): *Exists $\kappa > 1$ such that for all t, y_1, y_2 :*

$$\langle y_1 - y_2, s_\theta(t, y_1) - s_\theta(t, y_2) \rangle \leq -\kappa \|y_1 - y_2\|^2.$$

Theorem D.1. *If Assumption 4.1 and D.1 hold, $W_2(Y_T^{\text{SDE}}, Y_T^{\text{ODE}}) \leq \eta \|g\|_\infty \frac{\sqrt{A+2\kappa-2}}{\kappa-1}$.*

Proof. With the additional log-concavity assumption on the score function,

$$\mathbb{E}\langle Y_t^{\text{ODE}} - Y_t^{\text{SDE}}, s_\theta(t, Y_t^{\text{ODE}}) - s_\theta(t, Y_t^{\text{SDE}}) \rangle \leq -\kappa \cdot u(t).$$

In this case, the coefficient α becomes *contractive* as we may take $\delta = \frac{\kappa-1}{2\kappa}$:

$$\begin{aligned} u'(t) &\leq \|g\|_\infty^2 u(t) - \kappa \cdot \|g\|_\infty^2 \cdot u(t) + \eta^2 \|g\|_\infty^2 \left(\delta \cdot \kappa \cdot u(t) + \frac{A}{4\delta \cdot \kappa} \right) + \|g\|_\infty^2 \cdot \eta^2 \\ &\leq \underbrace{-\kappa \|g\|_\infty^2 \left(1 - \frac{1}{\kappa} - \delta \right)}_\alpha \cdot u(t) + \underbrace{\eta^2 \|g\|_\infty^2 \left(1 + \frac{A}{4\delta \cdot \kappa} \right)}_\beta. \end{aligned}$$

For arbitrary time horizon T , equation 18 gives:

$$\begin{aligned} u(T) &\leq \frac{(1 - e^{-\alpha T}) \cdot \beta}{-\alpha} \\ &\leq \frac{1}{\kappa \left(1 - \frac{1}{\kappa} - \delta \right)} \cdot \eta^2 \|g\|_\infty^2 \left(1 + \frac{A}{4\delta \cdot \kappa} \right) \\ &= \eta^2 \|g\|_\infty^2 \cdot \frac{2(\kappa - 1) + A}{(\kappa - 1)^2}. \end{aligned}$$

□

1242 *Remark.* Assumption 4.2 is valid for a *unimodal* terminal distribution, in which a strongly log-
1243 concave coefficient $\kappa(T)$ satisfies

$$1244 \mathbb{E}\langle Y_T^{\text{ODE}} - Y_T^{\text{SDE}}, s_\theta(T, Y_T^{\text{ODE}}) - s_\theta(T, Y_T^{\text{SDE}}) \rangle \leq -\kappa(T) \cdot u(T).$$

1245 In intermediate time steps,

$$1247 \mathbb{E}\langle Y_t^{\text{ODE}} - Y_t^{\text{SDE}}, s_\theta(t, Y_t^{\text{ODE}}) - s_\theta(t, Y_t^{\text{SDE}}) \rangle \leq -\kappa(t) \cdot u(t),$$

1248 in which (Tang & Zhao, 2025) bounds

$$1249 \kappa(t) \geq \frac{\kappa(T)}{\exp\left(-\int_0^{T-t} f(s) ds\right) + \kappa(T) \cdot \int_0^{T-t} \exp\left(-\int_s^{T-t} f(\tau) d\tau\right) ds}.$$

1252 Since $\kappa(t)$ is strictly positive on $[0, T]$, a global coefficient $\kappa_{\text{global}} = \inf_t \kappa(t)$ exists.

1254 D.3 DISCRETIZATION ERROR ANALYSIS

1255 Let N be the number of time steps and $h = T/N$ be the step size, we consider the DDPM sam-
1256 pling scheme with Euler–Maruyama discretization. Our Assumption 4.1 are equivalent to Assump-
1257 tions 6.1 and 6.2 (L -Lipschitz score and finite second moment) in (Liang et al., 2025), which devel-
1258 ops from (Chen et al., 2023a). For $h \lesssim 1/L$,

$$1260 \text{TV}(q_T^h, p_0) \lesssim \underbrace{\sqrt{\text{KL}(p_0 \|\mathcal{N}(0, I))e^{-T}}}_{\text{forward convergence}} + \underbrace{(L\sqrt{dh} + Lm_2h)\sqrt{T}}_{\text{discretization error}} + \underbrace{\varepsilon_0\sqrt{T}}_{\text{score error}},$$

1262 where q_T^h is the law of the discrete sampler at time T , d is the ambient space dimension, and $m_2^2 =$
1263 $\mathbb{E}_{p_0} \|X\|^2$. In particular, the discretization contribution to the total variation distance is $O(h^{1/2})$ as
1264 $h \rightarrow 0$ for fixed T, d, L, m_2 .

1266 Denote $Y_T^{\text{ODE}, h}$ and $Y_T^{\text{SDE}, h}$ as the ODE/SDE fine-tuned inference under step size h , by triangular
1267 inequality,

$$1268 W_2(Y_T^{\text{ODE}, h}, Y_T^{\text{SDE}, h}) \leq W_2(Y_T^{\text{ODE}}, Y_T^{\text{SDE}}) + \underbrace{W_2(Y_T^{\text{ODE}, h}, Y_T^{\text{ODE}}) + W_2(Y_T^{\text{SDE}}, Y_T^{\text{SDE}, h})}_{\varepsilon_d(h)}$$

1271 Since $\sup_t \mathbb{E} \|X_t\|^2 \leq \infty$, a standard coupling method between W_2 and TV distance gives.

$$1272 W_2(Y_T^\theta, Y_T^{\theta, h}) \lesssim \text{TV}(Y_T^\theta, Y_T^{\theta, h})^{1/2}.$$

1274 Therefore, the contribution of time discretization to the W_2 error is bounded by the additive term

$$1275 \varepsilon_d(h) := C (L\sqrt{dh} + Lm_2h)^{1/2},$$

1276 for some constant C depending only on T, d . In the continuous time limit, $\varepsilon_d(h) \rightarrow 0$.

1279 E DISCUSSION ON DDIM AND GDDIM UNDER HIGH STOCHASTICITY

1280 When $1.0 < \eta < +\infty$, the stochasticity level of gDDIM is always upper-bounded by

$$1282 \sigma_t^{\text{gDDIM}}(\eta) = (1 - \alpha_{t-\Delta t}) \left(1 - \left(\frac{1 - \alpha_{t-\Delta t}}{1 - \alpha_t} \right)^\eta \left(\frac{\alpha_t}{\alpha_{t-\Delta t}} \right)^\eta \right) \leq 1 - \alpha_{t-\Delta t} < +\infty;$$

1284 whereas the linear interpolation of DDIM gives an unbounded

$$1286 \sigma_t^{\text{DDIM}}(\eta) = \eta \sqrt{\frac{1 - \alpha_{t-\Delta t}}{1 - \alpha_t}} \sqrt{1 - \frac{\alpha_t}{\alpha_{t-\Delta t}}}.$$

1288 In order to ensure the well-posedness of $\sqrt{(1 - \alpha_{t-\Delta t} - \sigma_t^2)(1 - \alpha_t)}$ in equation 4, the conven-
1289 tional discretization requires

$$1290 \sigma_t^{\text{DDIM}} \leq 1 - \alpha_{t-\Delta t}.$$

1292 We may attempt to bypass the restriction by regulating

$$1293 \sigma_t^{\text{linear DDIM}}(\eta) = \min\{1 - \alpha_{t-\Delta t}, \eta_t^{\text{DDIM}}(\eta)\}.$$

1294 However, this changes the terminal marginals in equation 3. Therefore, the classical DDIM interpo-
1295 lation cannot simultaneously support arbitrary $\eta > 1.0$ and preserve the terminal marginal, whereas
gDDIM does.

1296 F HYPERPARAMETERS

1297
1298 All experiments are conducted on 7 Nvidia A100 GPUs. Mixed precision training is used with the
1299 bfloat16 (bf16) format.

1301 F.1 DDPO EXPERIMENTS

1302
1303 We follow the setup of Black et al. (2024), using denoising step $T = 50$ and guidance weight
1304 $w = 5.0$ throughout all experiments. We also use the AdamW optimizer Loshchilov & Hutter
1305 (2019) with default weight decay $1e-4$ and optimal learning rates for different reward functions.
1306 Reward gaps under four reward functions are shown in Figure 3 in Section 5 and Appendix H.

1307
1308 **Table 4:** DDPO hyperparameters

		ImageReward	PickScore	HPSv2	Aesthetic
	Batch size (Per-GPU)	48	24	24	32
DDPO	Samples per iteration (Global)	336	168	168	224
	Gradient updates per iteration	2	2	2	4
	Clip range	$1e-5$	$5e-5$	$1e-4$	$1e-4$
	Optimizer Learning Rate	$6e-4$	$6e-4$	$3e-4$	$3e-4$

1317
1318 **Animal Prompts** dataset consists of 398 animal labels extracted from *ImageNet-1k class labels*
1319 Deng et al. (2009), often with comma-separated synonyms and scientific names.

1320 **Comprehensive Prompts** dataset consists of 300 detailed and diverse descriptions of animals, ve-
1321 hicles, pieces of furniture, and landscapes with designated backgrounds, dynamics, or textiles.

1323 F.2 MIXGRPO EXPERIMENTS

1324
1325 We follow the setup of Li et al. (2025), letting reward model be "multi_reward" with equal weights.
1326 We set $T = 15$ as the denoising steps, AdamW optimizer with learning rate $1e-5$ and weight decay
1327 $1e-4$. For GRPO, the generation group size is 12 and clip range is $1e-4$. We perform 12 gradient
1328 updates per iteration.

1330 G LLM USAGE

1331
1332 Large Language Model (LLM) assists in LaTeX graphic alignments, spelling checks, and solving
1333 environment conflict issues in implementing DDPO and MixGRPO.

H MORE EXPERIMENT RESULTS

H.1 DDPO IMAGES ON DIFFERENT TRAINING STEPS



Figure 7: SDE (top) and ODE (bottom) sampling from every 100 training steps under PickScore with $\eta = 1.2$. Prompts (from left to right): “African chameleon, Chamae”, “Gordon setter”, “Great Pyrenees”, “malamute, malemute, Alask, Bra”, “Siamese cat, Siamese”, “bee eater”, “gaint schnauzer”, “Indian elephant, Elephas”, “marmoset”, “water buffalo, water ox”, “pug, pug dog”.

H.2 DDPO IMAGES UNDER IMAGEREWARD



Figure 8: Sampling schemes (from left to right): (i) SDE with $\eta = 0.75$, (ii) ODE with $\eta = 0.75$; (iii) SDE with $\eta = 1.2$, (iv) ODE with $\eta = 1.2$. Prompts (from top to bottom): “collie”, “old English sheepdog, bob”, “Irish terrier”.

H.3 DDPO IMAGES UNDER PICKSCORE



Figure 9: Sampling schemes (from left to right): (i) SDE with $\eta = 0.75$, (ii) ODE with $\eta = 0.75$; (iii) SDE with $\eta = 1.5$, (iv) ODE with $\eta = 1.5$. Prompts (from top to bottom): “baboon”, “white wolf, Arctic wolf”, “clumber, clumber spaniel”.

H.4 DDPO IMAGES UNDER MORE COMPREHENSIVE PROMPTS



Figure 10: ODE (below) image generation preserves prompt instructions with better quality on details compared to SDE (above) image generation under large stochasticity ($\eta = 1.2$).

Prompts (from left to right): “A vintage writing desk with an open journal and a flickering candle.”, “A macaque soaking in a steaming hot spring, surrounded by falling snow.”, “A wicker rocking chair on a wrap-around porch during golden hour.”, “A sleek sports car drifting on a mountain highway during golden hour.”

H.5 DDPO REWARD GAPS FOR OTHER REWARDS

H.5.1 IMAGEREWARD



Figure 11: Bounded reward gaps trained under ImageReward for 200 steps with stochasticity $\eta \in \{1.0, 1.2, 1.5\}$

1458 H.5.2 HPSv2
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481



1482 **Figure 12:** Bounded reward gaps trained under HPSv2 for 200 steps with stochasticity $\eta \in \{1.0, 1.2, 1.5\}$
 1483
 1484

1485 H.5.3 AESTHETIC
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511



Figure 13: Bounded reward gaps trained under Aesthetic for 200 steps with stochasticity $\eta \in \{1.0, 1.2, 1.5\}$