

INFERENCE-TIME DIFFUSION MODEL ALIGNMENT VIA RANDOM ORDINARY EQUATIONS



Figure 1: Our method successfully aligns diverse pre-trained diffusion models with various reward functions, such as aesthetics (a), semantics (d), compressibility (b), and sharpness and colors (e). It can also synergize with community-provided modules to improve the aesthetics of samples (c).

ABSTRACT

Aligning diffusion models (DM) with human preferences is a challenging yet practical task. Recent efforts focus on training-free methods, but usually adopt high-dimensional action spaces or require differentiable rewards. To address these issues, we propose a novel inference-time alignment framework based on random ordinary differential equation sampling. Specifically, we first formulate DM alignment as a max-encountered-reward optimal control problem. Then, by fixing the process noise and optimizing the perturbation strength, we obtain a 1-D action space, which integrates naturally with Monte Carlo tree search. We can thus perform trajectory search to derive the optimal control in a gradient-free manner, therefore supporting non-differentiable rewards. We also provide theoretical guarantees and empirical evidence to support and validate our method. Experiments show that our method demonstrates sufficient sample diversity and successfully aligns pre-trained DMs with reward functions defined on clean image domains. Our method outperforms traditional inference-step scaling, achieving higher best rewards. Meanwhile, it has significantly higher parameter efficiency than existing approaches adopting high-dimensional action spaces. Our approach can be plug-and-play integrated into any multi-step inference DMs.

1 INTRODUCTION

Diffusion models (DMs) (Ho et al., 2020; Song et al., 2021) can fit the data distribution, but often misaligns with human preferences, such as aesthetics (Fan et al., 2023), compressibility (Black et al., 2023), and colors (Eyring et al., 2024). To overcome this challenge, researchers (Lee et al., 2023; Black et al., 2023) model the denoising process as a Markov decision process (MDP), and train reward models to guide the training of DMs, thereby aligning them with specific rewards. But these approaches often require intensive training. A more favorable avenue is *inference-time scaling*, which employ inference-time guidance or optimization methods to achieve training-free DM alignment (Liu et al., 2024). However, these methods usually necessitate differentiable reward

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

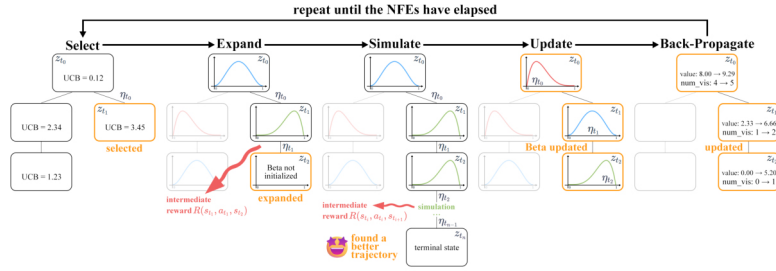


Figure 2: **Illustration of our MCTS.** The better trajectories found and the Beta policy updates in this figure occurred in the simulation phase. Remind that updates can also occur during the expansion phase.

functions (Bansal et al., 2023) and involve costly back-propagation (Eyring et al., 2024). Recently, researchers have turned to gradient-free methods to search noise or latents for inference-time alignment (Liu et al., 2024; Oshima et al., 2025), which, however, are always parameter-inefficient due to high-dimensional action spaces (Sec. 5.4). Besides, previous work also suffers from sparse rewards during the denoising process, resulting in low data efficiency (Zhang et al., 2024a).

To address these issues, we propose a novel inference-time scaling framework that alters the denoising trajectory to align DMs with any off-the-shelf reward function (unnecessarily differentiable). Specifically, our framework starts from a random ordinary differential equation (RODE) (Han et al., 2017) sampling extracted from DDIM (Song et al., 2020), which has lower variance and introduces controllable randomness to ODE sampling for exploration (Sec. 3.2). Then, we model the denoising process as a deterministic episodic MDP with a 1-D action space and dense rewards (Sec. 3.3), and formulate inference-time DM alignment as a max-encountered-reward optimal control problem (Quah & Quek, 2006) (Sec. 3.1, Sec. 3.4). Finally, we employ a value-based Monte Carlo tree search (MCTS) (Coulom, 2006; Browne et al., 2012) with online-updated Beta policies (Fig. 2) to derive the optimal control (Fleming & Rishel, 2012), thus achieving DM alignment (Sec. 4).

In experiments, we first verify that RODE sampling offers reasonable, sufficient and controllable diversity (Sec. 5.2). Then, we demonstrate that our combination of latent reward shaping, max-reward modeling and max-value update outperforms traditional ones (Sec. 5.3). Next, we compare different methods on aligning DMs to various reward functions, such as aesthetics (Sec. 5.4), semantics (Sec. 5.5) and compressibility (Sec. 5.6), highlighting the advantage of our MCTS with augmented RODE-based scaling. Finally, we present ablations and applications of our method (Sec. 5.7).

Our contributions: (i) We propose a novel parameter-efficient inference-time scaling framework for aligning DMs with any Lipschitz continuous reward functions (unnecessarily differentiable) by altering the denoising trajectory; (ii) **We are the first to** extract an RODE sampling from the DDIM’s SDE sampling, and model the denoising process as a MDP with a **1-D action space** and dense rewards; (iii) **We are the first to** formulate the inference-time DM alignment as a **max-encountered-reward optimal control problem**, and solve it using MCTS with augmented RODE-based scaling; Our approach can be plug-and-play integrated into any multi-step inference DMs.

2 PRELIMINARIES

2.1 DDIM

A stochastic DDIM (Song et al., 2020) step writes

$$z_{t_{i+1}} = \sqrt{\bar{\alpha}_{t_{i+1}}} \hat{z}_{t_n} + \sqrt{1 - \bar{\alpha}_{t_{i+1}} - \sigma_{t_i}^2 \epsilon_{\theta}^{(t_i)}}(z_{t_i}) + \sigma_{t_i} \epsilon_{t_i}, \quad (1)$$

where the process noise $\epsilon_{t_i} \sim \mathcal{N}(0, I)$ is re-sampled at each timestep t_i , contributing stochasticity via the perturbation term $\sigma_{t_i} \epsilon_{t_i}$; and the *posterior mean*

$$\hat{z}_{t_n} \triangleq \mathbb{E}[z_{t_n} | z_{t_i}] = (z_{t_i} - \sqrt{1 - \bar{\alpha}_{t_i}} \epsilon_{\theta}^{(t_i)}(z_{t_i})) / \sqrt{\bar{\alpha}_{t_i}}. \quad (2)$$

The standard deviation σ_{t_i} is defined as the product of η_{t_i} and a time-dependent constant ω_{t_i} :

$$\sigma_{t_i}(\eta_{t_i}) \triangleq \eta_{t_i} \omega_{t_i} \triangleq \eta_{t_i} \sqrt{(1 - \bar{\alpha}_{t_{i-1}}) / (1 - \bar{\alpha}_{t_i})} \sqrt{1 - \bar{\alpha}_{t_i} / \bar{\alpha}_{t_{i-1}}}. \quad (3)$$

η_{t_i} is conventionally stipulated to be time-invariant, *i.e.*, $\eta_{t_i} = \eta$ ($0 \leq i < n$). The value of η determines the stochasticity of the DDIM step: $\eta = 0$ makes the update process fully deterministic, referred to as *deterministic DDIM*; When $\eta \in (0, 1]$, the update process incorporates some randomness, termed *stochastic DDIM*. Particularly, it is called *DDPM* (Ho et al., 2020) when $\eta = 1$.

2.2 DDPO-STYLE MODELING

Most existing RL-based DM alignment methods (Black et al., 2023; Fan et al., 2023; Lee et al., 2023; Xu et al., 2023a; Clark et al., 2023; Prabhudesai et al., 2023) are based on the DDPO-style modeling (Black et al., 2023). Specifically, given a timestep schedule $\{t_0, t_1, \dots, t_n\}$, where $T = t_0 > t_1 > \dots > t_n = 0$, in which T represents the number of training steps of the DMs, and n denotes the number of inference steps. Meanwhile, given a DM p_θ , and a *clean reward model* $\phi(\cdot)$ (defined on clean image domains). Condition information (*e.g.*, text prompts) is omitted for simplicity. The denoising process of p_θ is modeled as a multi-step MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma, K)$, where:

$$s_{t_i} \triangleq (z_{t_i}, t_i), \quad \pi(a_{t_i}, s_{t_i}) \triangleq p_\theta(z_{t_{i+1}} | z_{t_i}, t_i), \quad a_{t_i} \triangleq z_{t_{i+1}}, \quad P(s_{t_{i+1}} | s_{t_i}, a_{t_i}) \triangleq (\delta_{z_{t_{i+1}}}, \delta_{t_{i+1}})$$

$$R(s_{t_i}, a_{t_i}, s_{t_{i+1}}) \triangleq \begin{cases} \mathbb{E}_{z_{t_n} \sim p_\theta(z_{t_{i+1}})}[\phi(z_{t_n})], & \text{if } i+1 = n \\ 0, & \text{otherwise} \end{cases}, \quad \gamma = 1, \quad K \triangleq n. \quad (4)$$

Refer to Appx. D.1 for symbol description and more discussion.

Given an initial noise z_{t_0} , the objective of RL is to maximize

$$J(\theta) \triangleq \mathbb{E}_{z_{t_n} \sim p_\theta(z_{t_0})} \left[\sum_{i=0}^{n-1} R(s_{t_i}, a_{t_i}) \right] \stackrel{\text{if DDPO}}{=} \mathbb{E}_{z_{t_n} \sim p_\theta(z_{t_0})} [\phi(z_{t_n})]. \quad (5)$$

2.3 LATENT REWARD SHAPING

To obtain dense rewards for the denoising process for faster convergence and better performance, training-free *latent reward shaping* utilizes the final reward (obtained with *clean reward models*) as a *proxy* for intermediate rewards (Appx. C.3). Specifically, to estimate the *latent reward* for an intermediate latent z_{t_i} ($0 \leq i < n$), denoted as $\hat{\phi}(z_{t_i})$, we first estimate the *pseudo-final sample* \hat{z}_{t_n} , from z_{t_i} in a certain way \mathcal{F} , *i.e.*, $\hat{z}_{t_n} = \mathcal{F}(z_{t_i})$. The estimated reward for z_{t_i} ($0 \leq i \leq n$) writes

$$\hat{\phi}(z_{t_i}) = \begin{cases} \phi(\mathcal{F}(z_{t_i})) = \phi(\hat{z}_{t_n}), & \text{if } i < n \\ \phi(z_{t_n}), & \text{if } i = n \end{cases}, \quad (6)$$

and the reward function in Eq. 4 is modified to

$$R(s_{t_i}, a_{t_i}, s_{t_{i+1}}) \triangleq \mathbb{E}_{z_{t_{i+1}} \sim p_\theta(\cdot | z_{t_i}, t_i)} [\hat{\phi}(z_{t_{i+1}})]. \quad (7)$$

For example, setting \mathcal{F} as the posterior mean of DDIM, *i.e.*, $\mathcal{F}(z_{t_i}) = \mathbb{E}[z_{t_n} | z_{t_i}]$ (Eq. 2), we can derive the *immediate-DDIM* paradigm latent reward shaping. Refer to Appx. D.2 for more details.

3 PROBLEM STATEMENT AND OUR FORMULATION

3.1 PROBLEM FORMULATION

Given a DM p_θ and a clean reward model $\phi(\cdot)$, the most straightforward problem formulation is to maximize the rewards of the final samples, *i.e.*, the *final reward* of the MDP, regardless of whether sparse or dense rewards are adopted. In the sparse reward scenario, the final reward is the only available object to optimize, which means the objectives of maximizing *cumulative rewards*, *final rewards*, and *rewards encountered during the process* are equivalent. However, we argue that, the above formulation is sub-optimal when adopting dense rewards (Zhang et al., 2024b;a) due to the inconsistency between intermediate rewards and final rewards, regardless of whether the final rewards are optimized directly (*e.g.*, intermediate rewards as heuristic) or indirectly (*e.g.*, maximizing

162 cumulative rewards), since intermediate rewards would mislead the search process due to similar
 163 scale and *inadmissibility* (Appx. E.1). To circumvent this issue, inspired by the *de novo* drug de-
 164 sign (Gummesson Svensson et al., 2024), we treat the pseudo-final samples obtained during the
 165 denoising process as valid samples, thereby transforming the problem into maximizing the encoun-
 166 tered rewards within an episode, *i.e.*, a *max-encountered-reward (max-reward for short) control/RL*
 167 *problem* (Quah & Quek, 2006; Gottipati et al., 2020) (Appx. D.3). **Note that this is different from**
 168 **purely maximizing the final reward, *i.e.*, *max-final-reward modeling*.** To the best of our knowledge,
 169 we are the first to model DM alignment as a max-reward control/RL problem.

170 Specifically, we adopt the DDIM schedule. Given an initial noise z_{t_0} and the NFE budget for
 171 calculating transition dynamics (*i.e.*, maximum inference steps) N , we aim to find the sample \tilde{z}_{t_N}
 172 within N DDIM steps that maximizes $\phi(\tilde{z}_{t_N})$, where \tilde{z}_{t_N} is defined as the derived sample z_{t_N} if
 173 N DDIM steps are strictly performed sequentially; otherwise, a pseudo-final sample $\tilde{z}_{t_N} = \mathcal{F}(z_{t_i})$.
 174 We do **not** enforce a full use of N steps, because: (i) Sabour et al. (2024); Ye et al. (2024);
 175 we found that, the fixed hand-crafted timestep schedule might be sub-optimal; (ii) Nichol & Dhariwal (2021);
 176 Li et al. (2023) showed that, increasing inference steps might cause sample quality degradation.

177 **Geometric Interpretation.** Fig. 3 provides
 178 a geometric interpretation of the conventional
 179 and our formulation. Under the conventional
 180 formulation, the MDP transfer graph G of the
 181 DM’s inference process is a directed tree due
 182 to strictly unidirectional and stepwise transi-
 183 tions. Instead, our formulation treats the in-
 184 process pseudo-final samples as valid samples,
 185 introducing *shortcuts* from states at timestep
 186 $t_i \in [t_1, t_{n-2}]$ to the terminal states at timestep
 187 t_n . This transforms G into a directed acyclic graph (DAG) G' (Bang-Jensen & Gutin, 2008), where
 188 nodes in each layer of G now have edges to both the next layer and the leaves.

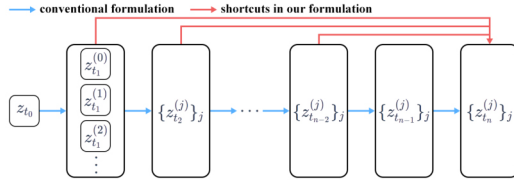


Figure 3: **Geometric interpretation of the conventional and our formulation.** Note that ours contains both blue and red edges.

189 **Enable Pseudo-Final Samples as Valid Samples.** We empirically observe reward hacking in
 190 our max-reward modeling due to early pseudo-final samples (Sec. 5.7). To mitigate this, we only
 191 consider pseudo-final samples with depth no less than $\tau \in [1, n]$ as valid samples. For samples with
 192 depth in the range $[1, \max\{\tau - 1, 1\}]$, we still compute intermediate rewards.

193
 194
 195 **3.2 AUGMENTED RODE SAMPLING FOR DDIM**

196 Existing inference-time scaling methods suffer from high-dimensional action spaces (Appx. C.2).
 197 To overcome this, we propose a random ODE-based (RODE-based) (Han et al., 2017) sampling
 198 for DDIM schedule to obtain a low-dimensional action space. It is directly extracted from the pre-
 199 trained SDE sampling without extra training or fine-tuning. Specifically, we fix a set of process
 200 noises $\{\epsilon_{t_i}\}_{i=0}^{n-1}$ for the DDIM step (Eq. 1), converting SDE sampling to ODE sampling, which is
 201 deterministic once given $\{\sigma_{t_i}\}_{i=0}^{n-1}$. Then, we propose to reintroduce randomness to the sampling
 202 process by treating the process variances $\{\eta_{t_i}\}_{i=0}^{n-1}$ in Eq. 3 as random variables, which are sampled
 203 from some certain probability density functions (PDFs). In the continuous-time perspective, this
 204 eliminates the $d\bar{w}$ term in the standard SDE sampling (Appx. D.4.1), transforming it into an RODE

205
 206
 207
$$\begin{aligned} \text{SDE: } dz &= [f(z, t) - g(t)^2 \nabla_z \log p_t(z)] dt + g(t) d\bar{w} \\ \text{RODE: } dz &= [f(z, t) - \omega(t)^2 \eta(t)^2 \nabla_z \log p_t(z)] dt + \omega(t) \eta(t) h(t) dt' \end{aligned} \tag{8}$$

208 where \bar{w} is an inverse Wiener process, $\eta(t) \in [0, 1]$ is a scalar stochastic process, and $h(t)$ is the
 209 fixed direction field induced by the realized process noise (ϵ_{t_k} in discrete representation).

210 Intuitively, our RDOE sampling fixes the perturbation direction at each timestep for all samples, and
 211 introduces randomness via random scaling along these directions. This yields lower variance than
 212 standard SDE sampling, where the process noise is resampled at each timestep:
 213

214 **Proposition 1.** *In the discrete analog, the variances of SDE and RODE step writes*

215
$$\text{Var}_{SDE} = \omega_{t_i}^2 \eta^2 \Delta t I, \quad \text{Var}_{RODE} = \omega_{t_i}^2 \text{Var}(\eta_{t_i}) (\Delta t)^2 \epsilon_{t_i} \epsilon_{t_i}^\top, \tag{9}$$

216 respectively. Thus, $\text{Var}_{\text{RODE}} < \text{Var}_{\text{SDE}}$ in Loewner order (Horn & Johnson, 2012), and both sampling
 217 schemes display exponential tail decay and strong concentration around their means.
 218

219 Despite lower variance, RODE sampling can provide sufficient sample diversity in multi-step infer-
 220 ence scenarios (Sec. 5.2). Besides, let $p^{(S)}$ and $p^{(R)}$ be the distribution induced by SDE and RODE
 221 sampling, respectively. $p^{(R)}$ approximates the true data distribution as $p^{(S)}$ does because:

222 **Proposition 2.** The Wasserstein-1 distance (Kolouri et al., 2017) $W_1(p^{(S)}, p^{(R)})$ is bounded:

223
 224
$$W_1(p^{(S)}, p^{(R)}) \leq \sum_{k=0}^{n-1} M_{t_{k+1}} (\omega_{t_k} \cdot C_{\epsilon, t_k} \cdot \mathbb{E}[\Delta\eta_{t_k}]) + \mathcal{O}\left(\sum_{k=0}^{n-1} M_{t_{k+1}} \gamma_{t_k}\right), \quad (10)$$

 225
 226

227 where: (i) $M_{t_{i+1}} = \prod_{j=i+1}^{n-1} L_{t_j}$ ($0 \leq j \leq n-1$) is the the product of a sequence of Lipschitz
 228 constants L_{t_j} with a uniform upper-bound $L < \infty$. We have $M_{t_{i+1}} \leq e^{LT}$, where T is the total time
 229 horizon (Gronwall, 1919); (ii) $C_{\epsilon, t_k} = \mathbb{E}[\|\epsilon_{t_k}\|]$; (iii) $\Delta\eta_{t_k} = |\eta_{t_k} - \eta| \leq 1$; (iv) the score-estimation
 230 error at timestep t_i is upper-bounded by a sub-Gaussian perturbation of scale γ_{t_i} .
 231

232 Hence, the discrepancy between samples generated from the two paradigms can be controlled, and
 233 the η -bias is the dominant source of discrepancy between RODE and SDE trajectories.

234 Additionally, reusing the SDE-trained score network $s_\theta(z_t, t) \approx \nabla_z \log p_t(z_t)$ within RODE sam-
 235 pling does **not** significantly increase score-estimation error since $s_\theta(\cdot)$ approximates marginal dis-
 236 tributions and is independent of sampling trajectories. Specifically,

237 **Proposition 3.** The difference of score-estimation error writes

238
 239
$$\left| \mathbb{E}_{p_t^{(R)}} \mathcal{E}(z, t) - \mathbb{E}_{p_t^{(S)}} \mathcal{E}(z, t) \right| \leq L_{\mathcal{E}, t} \cdot W_1(p_t^{(R)}, p_t^{(S)}), \quad (11)$$

 240

241 where $\mathcal{E}(z, t) = \|s_\theta(z, t) - s(z, t)\|$ is the point-wise score-estimation error at timestep t , and $s(z, t)$
 242 is the GT score function. Refer to Appx. F for assumptions, proofs and more discussion.
 243

244 Prop. 3 is empirically supported by the zero-shot FIDs (Yu et al., 2021; Podell et al., 2023) (Sec. 5.2).
 245 Fortunately, in our scenario, RODE sampling aims to introduce randomness for exploration, in-
 246 stead of exactly approximates the true data distribution. Besides, the theoretically larger error in
 247 RODE is visually imperceptible, *i.e.*, it is hard to distinguish RODE-sampled results from standard
 248 SDE/ODE-sampled ones (*e.g.*, Fig. 1). It indicates that, RODE sampling can still generate reason-
 249 able samples under larger error, which we attribute to the robustness of DMs and the drift term
 250 dominance (Appx. F.3). Still, we can design the process noise to span the ϵ_{t_i} space more effectively,
 251 enhancing exploration and reducing error compared to SDE sampling. This can be readily done by
 252 generating each process noise with distinct seeds, since high-dimensional random vectors are almost
 253 always nearly orthogonal to each other (Cai et al., 2013). Our RODE sampling introduces random-
 254 ness for exploration, while avoiding dealing with the stochastic optimal control problem (Evans,
 1983; Fleming & Rishel, 2012) and its notorious high variance (Sec. 3.4).
 255

256 **Stronger proposition.** Under stronger assumptions, we show mathematically that $p^{(R)}$ and $p^{(S)}$
 257 are sufficiently close (strengthened Prop. 2). Refer to Appx. F.3.3 for more discussion.
 258

259 **Benefits to the Search Process.** Standard stochastic DDIM sampling adopts $\eta_{t_i} = \eta$ for all $i \in$
 260 $[0, n-1]$, thereby anchoring a data manifold of a single noise-level $\sigma_{t_i}(\eta)$ for each timestep t_i
 261 (Eq. 3). In contrast, our method traverses multiple neighboring noise-level manifolds $\sigma_{t_i}(\eta_{t_i}^{(0)})$,
 262 $\sigma_{t_i}(\eta_{t_i}^{(1)})$, \dots by adapting different $\sigma_{t_i}(\eta_{t_i}^{(j)})$ on every visit to timestep t_i , rather than being confined
 263 to a single manifold as in Oshima et al. (2025).
 264

265 **Augmented RODE Sampling.** When only one set of process noise is realized, RODE sampling
 266 may under-explore directions that are under-represented (*e.g.*, directions orthogonal to all noise vec-
 267 tors), leading to sub-optimal rewards. To overcome this, we propose *augmented RODE sampling*
 268 (Aug. RODE sampling), which adopts several sets of process noise, and retains the optimal re-
 269 sult *per sample*. Formally, for an initial noise z_{t_0} , we take m independent sets of process noise
 $\{\{\epsilon_{t_i}^{(j)}\}_{i=0}^{n-1}\}_{j=0}^{m-1}$. We then search for each set and obtain the maximum rewards $\{r^{(j)}\}_{j=0}^{m-1}$, and

take $r^* = \max_j \{r^{(j)}\}_{j=0}^{m-1}$ as the maximum achievable reward for z_{t_0} . Intuitively, adopting more sets of process noise can better span the ϵ_{t_i} space, thus alleviating insufficient exploration.

3.3 MDP MODELING

Treating the η_{t_i} s as low-dimensional actions and their PDFs as policies forms a novel MDP modeling paradigm. We continue the symbols in Sec. 2.2 and Sec. 2.3, and only indicate the modified parts:

(b) The action space $\mathcal{A} = [0, 1]$, in which the policy and the action at time step t_i is modified to

$$\pi(a_{t_i}, s_{t_i}) \triangleq \tilde{\pi}(\eta_{t_i} \mid z_{t_i}, t_i), \quad a_{t_i} \triangleq \eta_{t_i}; \quad (12)$$

(c) The transition dynamics P is modified to

$$P(s_{t_{i+1}} \mid s_{t_i}, a_{t_i}) \triangleq p_\theta(z_{t_{i+1}} \mid z_{t_i}, \eta_{t_i}). \quad (13)$$

Although inherent stochastic in p_θ , the transition becomes deterministic when the process noise $\{\epsilon_{t_i}\}_{i=0}^{n-1}$ is fixed as done in RODE sampling. We further specify that $P(s_{t_{n+1}} \mid s_{t_n}, a) = \delta_{s_{t_n}}$ for any action a . We subsequently treat P as a deterministic function induced by a Lipschitz continuous model class (Asadi et al., 2018), and use $s_{t_{i+1}} = P(s_{t_i}, a_{t_i})$ for simplicity;

(d) The reward function is modified to

$$R(s_{t_i}, a_{t_i}, s_{t_{i+1}}) = \hat{\phi}(z_{t_{i+1}}). \quad (14)$$

The expectation is omitted due to the fully deterministic transition dynamics and the sole dependence on states (*i.e.*, trajectory-independent) for the rewards. We stipulate $R(s_{t_n}, a, \cdot) = 0$ for any action a . The reward function is assumed to be Lipschitz continuous;

(e) The discount factor γ is fixed to 1 because: (i) we treat pseudo-final samples as valid ones in the max-reward modeling, thus no discount is needed; (ii) the episodic MDP is finite, so the Bellman update (Appx. D.3) converges without requiring $\gamma < 1$. Refer to Appx. E.2 for more discussion.

3.4 OPTIMAL CONTROL

We formulate the training-free inference-time DM alignment as an online planning problem (Efroni et al., 2020). Specifically, let $\eta_{t_0, \dots, t_{n-1}} = \{\eta_{t_0}, \dots, \eta_{t_{n-1}}\}$, and the optimization object be $J(\eta_{t_0, \dots, t_{n-1}})$. Our goal is to find an optimal control sequence $\eta_{t_0, \dots, t_{n-1}}^*$ that maximizes J , *i.e.*,

$$\eta_{t_0, \dots, t_{n-1}}^* = \arg \max_{\eta_{t_0, \dots, t_{n-1}}} J(\eta_{t_0, \dots, t_{n-1}}) \text{ s.t. } s_{t_{i+1}} = P(s_{t_i}, a_{t_i}). \quad (15)$$

We keep the DM’s parameters unchanged during the optimization process (*i.e.*, as a fixed environment model), and only optimize the policies used to sample η_t s, which prevents fatal deviations from the data distribution and unreasonable over-optimization for higher rewards (Black et al., 2023; Zhang et al., 2024b). Unlike He et al. (2025), our modeling does **not** require an additional KL divergence term for regularization. In summary, our RODE-based scaling adopts a low-dimensional action space that can be efficiently covered, enabling controllable trajectory modulation and smooth navigation of the continuous action space. Refer to Appx. E.3 and Appx. E.4 for more discussion.

4 ONLINE PLANNING USING MCTS

We employ MCTS as a training-free optimal control solver to derive the $\eta_{t_0, \dots, t_{n-1}}^*$. The overview is presented in Fig. 2. Specifically, to handle the continuous action space, we select node with the highest *value-based* UCB (Sec. 4.1), then expand it with Beta policies (Sec. 4.2). In the expansion and simulation phases, the *unimodal Beta policies* of the tree nodes are updated online if superior trajectories are found. Finally, the roll-out is backpropagated with *max-value policies* (Sec. 4.1).

4.1 SELECTION POLICIES

Value-Based UCB. We propose a value-based UCB to facilitate the calculation of UCB in continuous action space based on the insight that, the state value function $V(s)$ aggregates the performance

of all sampled actions at state s . Specifically, we propose to replace the $Q(s, a)$ in traditional UCB (Appx. D.5) with $V(s)$, and modify the exploration term in terms of s . Formally,

$$\text{exploitation}(s) = V(s), \quad \text{exploration}(s) = \sqrt{\ln N(\text{parent}(s))/N(s)}, \quad (16)$$

where $\text{parent}(s)$ and $N(s)$ denote the parent node and the visit times of state s , respectively.

Max-Value Policies. In cumulative-reward RL, the $V(s)$ is typically estimated with the average return from all roll-outs passing through s , which is updated with the *average-value policy* during the backpropagation phase (Appx. D.5). However, our max-reward modeling leads us to a novel *max-value policy*, where the $V(s)$ is defined and updated as

$$V(s) = \max_{i=0}^{N(s)-1} R_i, \quad V(s) \leftarrow \max\{V(s), R\}, \quad (17)$$

in which R is the cumulative or maximum return from the expansion or simulation phases.

Selection Depth Limit. During the MCTS selection phase, the depth of the node selected for expansion is limited to the range $[0, n']$, where $n' \leq n$. This is based on the observation from [Cao et al. \(2023\)](#); [Yang et al. \(2023\)](#) that, DMs recover the layout of images first, and the details later. We allocate more computational resources to shallower nodes in MCTS, which focuses on early-to-mid denoising steps, enabling more significant changes for potentially higher rewards.

4.2 BETA POLICY FOR CONTINUOUS ACTION SPACE

Unimodal Beta Policies. For continuous action spaces, uniform random sampling (Sec. 5.7) or discretization ([Chou et al., 2017](#)) can lead to sub-optimal performance. Following [Ye et al. \(2024\)](#), we maintain a unimodal (for policy stability) Beta distribution ([Gupta & Nadarajah, 2004](#)) for each MCTS node, and sample actions η_t from it, *i.e.*, we adopt Beta policies. To do this, we re-parameterize the shape parameters α and β (**not** the variance parameters α_{t_i} and β_{t_i} in DMs) of the Beta distribution as $\alpha = 1 + e^a$, $\beta = 1 + e^b$, and maintain the $a, b > 0$ for each MCTS node.

Mode Re-Parameterization. The Beta policies should be updated online during the search process, allocating more computational resources to promising actions. Directly update the shape parameters is hard to balance the exploration and exploitation of the Beta policy (Appx. G.2). To overcome this, we note that, the goal of the update of Beta policies is to maximize the likelihood of sampling (near-)optimal actions, and the main probability density of a unimodal Beta distribution is concentrated near its mode. This leads us to update with the mode. Specifically, we re-parameterize the Beta distribution in terms of its mode, and introduce a hyper-parameter $\zeta > 0$ to control its concentration. Specifically, let $\rho = (\alpha - 1)/(\alpha + \beta - 2)$ be the mode, then the a, b are re-parameterized as $a = \ln \zeta$, $b = \ln \zeta + \ln((1 - \rho)/\rho)$. Refer to Appx. G.2 for more discussion.

Initialization and Online Update. The policies for sampling actions are initialized to uniform distributions at the beginning, and then updated to unimodal Beta distributions with the first-sampled action. When a trajectory superior to the best-known one is found, the Beta policies from the current node to the root node in the MCTS tree are softly updated (Appx. G.3). Then, the a and b are updated correspondingly. In MCTS, such updates may occur within both the expansion (when expanding to a terminal state) and simulation phases, *i.e.*, we also treat the trajectories obtained via simulation as valid ones for updating Beta policies. Particularly, in our max-reward modeling, updates might also occur during the expansion phase — updated with pseudo-final samples when calculating the intermediate rewards for the expanded nodes. Regardless of where the update occurs, only the tree portion of the trajectory (real tree nodes) will have their Beta policies updated.

5 EXPERIMENTS

5.1 SETTINGS AND EVALUATIONS

Settings. We adopt the following settings unless specified. We generate $N = 2$ images per prompt using N fixed independent initial noise. We fix $m = 3$ sets of process noise for Aug. RODE sampling. For SDE sampling, DDPM adopts $\eta_{t_i} = 1.0$ for all $i \in [0, n - 1]$, and is run with $m = 3$

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

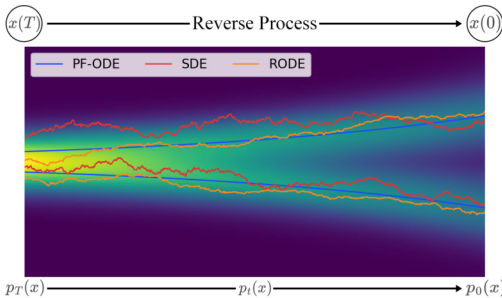


Figure 4: Visualization of denoising trajectories of 1× DDPM, 1× DDIM, and 2× ours.

Table 1: PF-ODE/SDE/RODE trajectories.

Table 2: Policies and MPDs (↑) of different inference-time scaling paradigms. Cells with a red / orange / yellow background: the best / second-best / third-best performance for each column.

Paradigm	Sampler	Components			SD-Turbo		SD v1.4	
		$\{z_{t_0}^{(j)}\}_{j=0}^{N-1}$	$\{\epsilon_{t_i}\}_{i=0}^{n-1}$	$\{\eta_{t_i}\}_{i=0}^{n-1}$	2-step	15-step	20-step	25-step
a) initial-noise	ODE	random	×	$\eta_{t_i} = 0$	0.5867	0.7542	0.7584	0.7594
	SDE	random	random	$\eta_{t_i} = 1.0$	0.6437	0.7653	0.7617	0.7650
	Aug. RODE	random	m sets	$\eta_{t_i} = 1.0$	0.6200	0.6861	0.6887	0.6914
b) process-noise	SDE	fixed	random	$\eta_{t_i} = 1.0$	0.3823	0.6848	0.6939	0.6946
c) process-noise -variance	Aug. RODE	fixed	m sets	$\eta_{t_i} \sim \mathcal{U}([0, 1])$	0.1628	0.4717	0.4852	0.4899
	Aug. RODE	fixed	m sets	$\eta_{t_i} \sim \mathcal{U}([0.5, 1])$	0.1295	0.3556	0.3649	0.3686

random sets of process noise, and retain the best results across them. In our max-reward formulation, $\tau = \lceil n/3 \rceil$. For MCTS, $n' = \lceil (4n)/5 \rceil$, and simulation actions are uniformly sampled from $[0, 1]$. For Beta policies, $\zeta = 3$. We adopt *immediate-ddim* policy, and the combination of latent reward shaping, max-reward modeling and max-value update for its best performance (Sec. 5.3).

Evaluations. We adopt “Best Reward” (the highest reward achieved), “Parameter Efficiency” and “Time Cost” (the averaged wall-clock time required per sample to obtain Best Reward) for evaluation. Refer to Appx. H for details on settings, evaluations, and experimental composition.

Table 3: FIDs (↓) for different sampling paradigms.

Paradigm / Step	15	20	25	30	50
ODE	29.4050	29.3332	29.4872	29.5351	29.7702
SDE	25.0427	24.9356	24.9540	24.6273	24.6893
Aug. RODE	28.0226	27.2009	27.1311	27.0817	26.9312

5.2 RODE SAMPLING

Trajectory Simulation. Fig. 1 visualizes the sampling trajectories produced by three paradigms. PF-ODE (Song et al., 2021) yields smooth trajectories, whereas both SDE and RODE generate jagged paths, where the RODE trajectories exhibits markedly smaller oscillations than the SDE ones, empirically supporting Prop. 1. Refer to Appx. I.1 for details.

Empirical Distribution. We compute zero-shot FID (Podell et al., 2023) on the MS-COCO 5k validation set (Lin et al., 2014) for different sampling paradigms. Results in Tab. 3 demonstrate that, our novel paradigm does not excessively deviate from the real data distribution, empirically corroborating Prop. 2. Refer to Appx. I.3 for more details.

Sample Diversity. We generate images under the conditions in Tab. 2, and compute the mean pairwise distance (MPD ↑) (Appx. H.3.1) to evaluate the sample diversity. It demonstrates that: (i) SDE sampling sustains the highest diversity, while our Aug. RODE sampling exhibits lower sample diversity than SDE sampling, which empirically validated Prop. 1; (ii) Our RODE sampling provides mid-to-high and controllable diversity when n is large, suggesting that, our method is more suitable for multi-step inference scenarios due to broader exploration. Refer to Appx. I.2 for more details.

Trajectory Visualization. The sampling trajectory of RODE sampling can be efficiently visualized with $\eta_{t_0, \dots, t_{n-1}}$ (Appx. I.4). Fig. 4 shows that, RODE sampling explores various paths from the origin to the n -th level data manifold (the outmost semi-circle), seeking for higher rewards.

Table 4: Comparison between different reward shaping, MDP modeling and value update policies. “latent”: latent reward shaping; “sparse”: w/o reward shaping. “cum.”: cumulative-reward modeling. “aver.”: average-value update.

Reward Shaping	MDP Modeling	Value Update	Best Reward
latent	cum.	aver.	28.9145
latent	cum.	max	28.8445
latent	max	aver.	28.8480
latent	max	max	29.1152
sparse	cum.	aver.	28.6840
sparse	cum.	max	29.0621

Table 5: Comparison between different methods when optimizing HPS v2. “-eps” / “-eta”: ϵ_{t_i} s / η_{t_i} s as actions. Note that Best-of-N fall outside our setting.

Method	Best Reward	Relative Impro. (%)	Param. Efficiency	Time Cost
DDIM (15 steps, baseline)	26.2122	/	/	147
Best-of-N (DDIM)	28.4056	8.37	2.03e-6	4480
Best-of-N (DDPM)	28.4297	8.46	2.05e-6	4513
DDIM (999 steps)	26.5362	1.24	/	6694
DDPM (200 steps)	27.6256	5.39	1.4378e-7	1382
GS-eps	29.7386	13.45	4.7830e-6	5684
BS-eps	30.2095	15.25	5.4217e-6	5469
MCTS-eps	29.0241	10.73	3.8139e-6	1437
GS-eta	27.8116	6.10	0.0355	5697
BS-eta	28.0833	7.14	0.0416	5273
Ours	28.1644	7.45	0.0434	2736

5.3 COMPARISON BETWEEN REWARD SHAPING, MDP MODELING AND VALUE POLICIES

This section compares various combinations of reward shaping, MDP modeling, and value update policies. Specifically, we adopt uniform expansion policy for all combinations, and align 15-step SD v1.4 (Rombach et al., 2022) to HPS v2 (Wu et al., 2023a). Results in Tab. 4 highlights the advantage of adopting latent reward shaping, max-reward modeling and max-value update.

5.4 ALIGNING WITH AESTHETICS

We align 15-step SD v1.4 (Rombach et al., 2022) to HPS v2 (Wu et al., 2023a), and make comparison between: (i) vanilla 15-step DDIM as baseline; (ii) best-of-N (Ma et al., 2025) with $N = 66$, which is not directly comparable as it falls outside our setting, but is included here as a widely used baseline; (iii) the inference-step scaling paradigm (Nichol & Dhariwal, 2021); (iv) beam search (BS) (Oshima et al., 2025) with $\epsilon_{t_i}/\eta_{t_i}$ actions; (v) greedy search (GS, the BS with a single beam); (vi) our MCTS. Tab. 5 demonstrates that: (i) Ours (MCTS-eta) surpasses traditional inference-step scaling. Note that scaling DDPM is essentially an ϵ_{t_i} -action method; (ii) When scaling with η_{t_i} , ours outperform both GS and BS, highlighting MCTS’s ability to allocate NFEs judiciously, and Beta policies’ capacity for accurate low-dimensional search; (iii) Although ϵ_{t_i} -action methods achieve the highest absolute score, its relative improvement is only $\sim 2\times$ better than η_{t_i} -action ones, but with $\sim 16k\times$ more parameters, revealing severe parameter redundancy; (iv) When naively transferring our MCTS to ϵ_{t_i} -actions, performance drops below GS and BS, underscoring the limitation of vanilla MCTS in high-dimensional continuous action spaces (Bianchi et al., 2023). Qualitative results are shown in Fig. 1 (a), where trajectory search-based methods yield samples that are visually more appealing than those produced by inference-step scaling. Refer to Appx. J for more details.

Besides, our method can generalize to other latent reward policies, see Appx. M.

5.5 ALIGNING WITH SEMANTICS

We align 30-step SDXL base (Podell et al., 2023) to CLIP score (Hessel et al., 2021), and benchmark our method against DDPM, DDIM and Z-Sampling (Bai et al., 2024). Qualitative results in Fig. 1 (d) are generated with text prompts “A giraffe underneath a microwave.” and “A red colored banana.”, respectively. It can be observed that, our method renders accurate positional relationship (line 1) and colors (line 2), while baselines fail to capture. Quantitative results are shown in Tab. 6. Refer to Appx. K for more details.

Table 6: Comparison between different methods when optimizing CLIP Score.

Method	Ours	DDPM	DDIM	Z-Sampling
CLIP score	0.3716	0.3569	0.3328	0.3676
Improvement (%)	11.66	7.24	/	10.46
HPS v2	29.0685	28.6943	27.8625	29.5635
Improvement (%)	4.33	2.99	/	6.10

Table 7: Comparison between different optimization objectives. “Ours (R)”: aligning with reward R with our method.

Method	CR (\uparrow)	HPS v2 (\uparrow)	CLIP Score (\uparrow)
Ours (CR)	2.9308	26.7430	0.3404
Ours (HPS v2)	2.9261	27.8809	0.3731
Ours (Composite)	2.9282	27.7684	0.3743

5.6 ALIGNING WITH COMPOSITE REWARDS

We align 50-step Pixart- α (Chen et al., 2023b) (DiT-architecture (Peebles & Xie, 2023)) to compressibility reward (CR) (Appx. H.3.4). Samples displayed in Fig. 1 (b) are generated with the text prompt “A beautiful blue and pink sky overlooking the beach.”, which introduces unreasonable smooth to inflate the reward. To overcome this, we optimize a *composite reward* that combines CR and HPS v2. Qualitative results demonstrate that, optimizing the composite reward yields samples that are perceptually smoother yet retain clear semantics. Quantitative results in Tab. 7 also highlights its balanced performance than optimizing either reward in isolation. Refer to Appx. L for more details.

5.7 ABLATIONS AND APPLICATIONS

Reward Hacking and Effects of τ . We align 50-step SD v1.4 to Laplacian variance (LAPV) and color channel reward (CCR) (Appx. H.3.4) with $\tau = 17$ (default) and $\tau = 1$ (enabling pseudo-final samples throughout the entire process), respectively. Fig. 1 (e) illustrates the reward hacking phenomena that leverage the noisy and blurred pseudo-final samples in the early denoising for higher rewards. This can be mitigated by adjusting τ to determine when to enable treating pseudo-final samples as valid ones. Refer to Appx. N.1 for more details. In summary, our method generalizes across diverse DM architectures, sampling conditions, inference steps, reward functions, *et al.*

Main Ablations. We conduct ablations on key components of our method. Results in Tab. 8 underscore advantage of our max-encountered-reward formulation over the conventional max-final-reward formulation. Refer to Appx. N and Appx. P for more ablations and failure cases, respectively.

Applications. Our method can synergize with Promptist (Hao et al., 2023) and Golden Noise (GN) (Zhou et al., 2024b) to improve sample aesthetics (Fig. 1 (c), Appx. O.1). Additionally, it can be used to quantitatively assess the robustness of reward models (Appx. O.2).

6 RELATED AND CONCLUSION

Related. Appx. C presents literature review on DM alignment, inference-time scaling for DMs and latent reward shaping. Our method belongs to training-free alignment via trajectory search.

Conclusion. This study makes the first step towards scaling DMs with 1-D actions. We present an augmented RODE-based scaling paradigm, and achieve parameter-efficient inference-time DM alignment by solving a max-reward optimal control problem using MCTS. Our method do **not** necessitate differentiable reward functions, and is particularly suitable for multi-step inference scenarios. Refer to Appx. A for use of LLMs, and Appx. Q for broader impacts and limitations.

REPRODUCIBILITY STATEMENT

Refer to Appx. H for general experimental settings. Detail settings for Sec. 5 are detailed in Appx. I to Appx. P. Refer to *supplementary materials* for code implementation.

Table 8: Ablations. “w/o pseudo-final”: disabling pseudo-final samples as valid samples. “w/o Beta policy”: uniform policy.

Method	Best Reward
Ours	29.2582
Ours w/o pseudo-final	29.0199
Ours w/o expansion depth limit	29.1082
Ours w/o online update	29.1547
Ours w/o Beta policy	29.1152

REFERENCES

- 540
541
542 Donghoon Ahn, Jiwon Kang, Sanghyun Lee, Jaewon Min, Minjae Kim, Wooseok Jang, Hyoungwon
543 Cho, Sayak Paul, SeonHwa Kim, Eunju Cha, et al. A noise is worth diffusion guidance. *arXiv*
544 *preprint arXiv:2412.03895*, 2024.
- 545
546 Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic inter-
547 polants. *arXiv preprint arXiv:2209.15571*, 2022.
- 548
549 Kavosh Asadi, Dipendra Misra, and Michael Littman. Lipschitz continuity in model-based rein-
550 forcement learning. In *International Conference on Machine Learning*, pp. 264–273. PMLR,
551 2018.
- 552
553 Lichen Bai, Shitong Shao, Zikai Zhou, Zipeng Qi, Zhiqiang Xu, Haoyi Xiong, and Zeke Xie.
554 Zigzag diffusion sampling: Diffusion models can self-improve via self-reflection. *arXiv preprint*
555 *arXiv:2412.10891*, 2024.
- 556
557 Jørgen Bang-Jensen and Gregory Z Gutin. *Digraphs: theory, algorithms and applications*. Springer
558 Science & Business Media, 2008.
- 559
560 Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas
561 Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the*
562 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 843–852, 2023.
- 563
564 Shane Barratt and Rishi Sharma. A note on the inception score. *arXiv preprint arXiv:1801.01973*,
565 2018.
- 566
567 Farhat Lamia Barsha and William Eberle. An in-depth review and analysis of mode collapse in
568 generative adversarial networks. *Machine Learning*, 114(6):141, 2025.
- 569
570 Mihaly Bencze, Constantin P Niculescu, and Florin Popovici. Popoviciu’s inequality for functions
571 of several variables. *Journal of mathematical analysis and applications*, 365(1):399–409, 2010.
- 572
573 Federico Bianchi, Lorenzo Bonanni, Alberto Castellini, Alessandro Farinelli, et al. Monte carlo tree
574 search planning for continuous action and state space. In *AIRO 2022 Artificial Intelligence and*
575 *Robotics 2022*, pp. 38–47. 2023.
- 576
577 Jeremiah Birrell. Concentration inequalities for the stochastic optimization of unbounded objectives
578 with application to denoising score matching. *arXiv preprint arXiv:2502.08628*, 2025.
- 579
580 Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion
581 models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- 582
583 Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp
584 Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey
585 of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in*
586 *games*, 4(1):1–43, 2012.
- 587
588 Tony Cai, Jianqing Fan, and Tiefeng Jiang. Distributions of angles in random packing on spheres.
589 *The Journal of Machine Learning Research*, 14(1):1837–1864, 2013.
- 590
591 Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Mas-
592 actrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. In
593 *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 22560–22570,
2023.
- Defang Chen, Zhenyu Zhou, Jian-Ping Mei, Chunhua Shen, Chun Chen, and Can Wang. A geomet-
ric perspective on diffusion models. *arXiv preprint arXiv:2305.19947*, 2023a.
- Defang Chen, Zhenyu Zhou, Can Wang, Chunhua Shen, and Siwei Lyu. On the trajectory regularity
of ode-based diffusion sampling. *arXiv preprint arXiv:2405.11326*, 2024.
- Defang Chen, Zhenyu Zhou, Can Wang, and Siwei Lyu. Geometric regularity in deterministic
sampling of diffusion-based generative models. *arXiv preprint arXiv:2506.10177*, 2025a.

- 594 Huanran Chen, Yinpeng Dong, Zeming Wei, Hang Su, and Jun Zhu. Towards the worst-case robust-
595 ness of large language models. *arXiv preprint arXiv:2501.19040*, 2025b.
596
- 597 Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James
598 Kwok, Ping Luo, Huchuan Lu, et al. Pixart-alpha: Fast training of diffusion transformer for
599 photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023b.
- 600 Po-Wei Chou, Daniel Maturana, and Sebastian Scherer. Improving stochastic policy gradients in
601 continuous control with deep reinforcement learning using the beta distribution. In *International
602 conference on machine learning*, pp. 834–843. PMLR, 2017.
603
- 604 Hyungjin Chung, Jeongsol Kim, Michael T Mccann, Marc L Klasky, and Jong Chul Ye. Diffusion
605 posterior sampling for general noisy inverse problems. *arXiv preprint arXiv:2209.14687*, 2022.
- 606 Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models
607 on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.
608
- 609 Adrien Couëtoux, Jean-Baptiste Hoock, Nataliya Sokolovska, Olivier Teytaud, and Nicolas Bon-
610 nard. Continuous upper confidence trees. In *International conference on learning and intelligent
611 optimization*, pp. 433–445. Springer, 2011.
- 612 Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *International
613 conference on computers and games*, pp. 72–83. Springer, 2006.
614
- 615 Rémi Coulom. Computing “elo ratings” of move patterns in the game of go. *ICGA journal*, 30(4):
616 198–208, 2007.
- 617 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances
618 in neural information processing systems*, 34:8780–8794, 2021.
619
- 620 Yonathan Efroni, Mohammad Ghavamzadeh, and Shie Mannor. Online planning with lookahead
621 policies. *Advances in Neural Information Processing Systems*, 33:14024–14033, 2020.
622
- 623 Jacob Eisenstein, Chirag Nagpal, Alekh Agarwal, Ahmad Beirami, Alex D’Amour, DJ Dvijotham,
624 Adam Fisch, Katherine Heller, Stephen Pfohl, Deepak Ramachandran, et al. Helping or herd-
625 ing? reward model ensembles mitigate but do not eliminate reward hacking. *arXiv preprint
626 arXiv:2312.09244*, 2023.
- 627 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
628 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers
629 for high-resolution image synthesis. In *Forty-first international conference on machine learning*,
630 2024.
- 631 Lawrence C Evans. An introduction to mathematical optimal control theory version 0.2. *Lecture
632 notes available at http://math.berkeley.edu/evans/control_course.pdf*, 1983.
633
- 634 Luca Eyring, Shyamgopal Karthik, Karsten Roth, Alexey Dosovitskiy, and Zeynep Akata. Reno:
635 Enhancing one-step text-to-image models through reward-based noise optimization. *Advances in
636 Neural Information Processing Systems*, 37:125487–125519, 2024.
- 637 Ying Fan and Kangwook Lee. Optimizing ddpm sampling with shortcut fine-tuning. *arXiv preprint
638 arXiv:2301.13362*, 2023.
639
- 640 Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,
641 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for
642 fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*,
643 36:79858–79885, 2023.
- 644 Charles Fefferman, Sanjoy Mitter, and Hariharan Narayanan. Testing the manifold hypothesis.
645 *Journal of the American Mathematical Society*, 29(4):983–1049, 2016.
646
- 647 Guilherme Fernandes, Vasco Ramos, Regev Cohen, Idan Szpektor, and João Magalhães. Latent
beam diffusion models for decoding image sequences. *arXiv preprint arXiv:2503.20429*, 2025.

- 648 Wendell H Fleming and Raymond W Rishel. *Deterministic and stochastic optimal control*, volume 1. Springer Science & Business Media, 2012.
- 649
- 650
- 651 Christine Fricker, Philippe Robert, and James Roberts. A versatile and accurate approximation for lru cache performance. In *2012 24th international teletraffic congress (ITC 24)*, pp. 1–8. IEEE, 2012.
- 652
- 653
- 654 GeeksforGeeks. How to check for blurry images in your dataset using the laplacian method, 2024. URL <https://www.geeksforgeeks.org/how-to-check-for-blurry-images-in-your-dataset-using-the-laplacian-method/>.
- 655
- 656
- 657
- 658 Andrew S Glassner. *An introduction to ray tracing*. Morgan Kaufmann, 1989.
- 659
- 660 Sai Krishna Gottipati, Yashaswi Pathak, Rohan Nuttall, Raviteja Chunduru, Ahmed Touati, Sri-ram Ganapathi Subramanian, Matthew E Taylor, Sarath Chandar, et al. Maximum reward formulation in reinforcement learning. *arXiv preprint arXiv:2010.03744*, 2020.
- 661
- 662
- 663 Thomas Hakon Gronwall. Note on the derivatives with respect to a parameter of the solutions of a system of differential equations. *Annals of Mathematics*, 20(4):292–296, 1919.
- 664
- 665
- 666 Marek Grzes. Reward shaping in episodic reinforcement learning. 2017.
- 667
- 668 Hampus Gummesson Svensson, Christian Tyrchan, Ola Engkvist, and Morteza Haghiri Chehreghani. Utilizing reinforcement learning for de novo drug design. *Machine Learning*, 113(7):4811–4843, 2024.
- 669
- 670
- 671 Xiefan Guo, Jinlin Liu, Miaomiao Cui, Jiankai Li, Hongyu Yang, and Di Huang. Initno: Boosting text-to-image diffusion models via initial noise optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9380–9389, 2024.
- 672
- 673
- 674 Arjun K Gupta and Saralees Nadarajah. *Handbook of beta distribution and its applications*. CRC press, 2004.
- 675
- 676
- 677 Xiaoying Han, Peter E Kloeden, Xiaoying Han, and Peter E Kloeden. *Random ordinary differential equations*. Springer, 2017.
- 678
- 679 Yaru Hao, Zewen Chi, Li Dong, and Furu Wei. Optimizing prompts for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:66923–66939, 2023.
- 680
- 681
- 682 Joshua Hare. Dealing with sparse rewards in reinforcement learning. *arXiv preprint arXiv:1910.09281*, 2019.
- 683
- 684 Haoran He, Jiajun Liang, Xintao Wang, Pengfei Wan, Di Zhang, Kun Gai, and Ling Pan. Scaling image and video generation via test-time evolutionary search. *arXiv preprint arXiv:2505.17618*, 2025.
- 685
- 686
- 687
- 688 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- 689
- 690 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- 691
- 692
- 693 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- 694
- 695
- 696 Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- 697
- 698 Robert William Gainer Hunt. *The reproduction of colour*. John Wiley & Sons, 2005.
- 699
- 700 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL <https://doi.org/10.5281/zenodo.5143773>. If you use this software, please cite it as below.
- 701

- 702 Vineet Jain, Kusha Sareen, Mohammad Pedramfar, and Siamak Ravanbakhsh. Diffusion tree sam-
703 pling: Scalable inference-time alignment of diffusion models. *arXiv preprint arXiv:2506.20701*,
704 2025.
- 705 Ioannis Karatzas and Steven Shreve. *Brownian motion and stochastic calculus*, volume 113.
706 Springer Science & Business Media, 2012.
- 707 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the de-
708 sign space of diffusion-based generative models. In S. Koyejo, S. Mohamed,
709 A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Infor-*
710 *mation Processing Systems*, volume 35, pp. 26565–26577. Curran Associates, Inc.,
711 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/](https://proceedings.neurips.cc/paper_files/paper/2022/file/a98846e9d9cc01cfb87eb694d946ce6b-Paper-Conference.pdf)
712 [file/a98846e9d9cc01cfb87eb694d946ce6b-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/a98846e9d9cc01cfb87eb694d946ce6b-Paper-Conference.pdf).
- 713 Sunwoo Kim, Minkyu Kim, and Dongmin Park. Alignment without over-optimization: Training-
714 free solution for diffusion models. *arXiv preprint arXiv:2501.05803*, 2025.
- 715 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Ad-*
716 *vances in neural information processing systems*, 34:21696–21707, 2021.
- 717 Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013.
- 718 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-
719 a-pic: An open dataset of user preferences for text-to-image generation. *Advances in Neural*
720 *Information Processing Systems*, 36:36652–36663, 2023.
- 721 Soheil Kolouri, Se Rim Park, Matthew Thorpe, Dejan Slepcev, and Gustavo K Rohde. Optimal
722 mass transport: Signal processing and machine-learning applications. *IEEE signal processing*
723 *magazine*, 34(4):43–59, 2017.
- 724 Lauwerens Kuipers and Harald Niederreiter. *Uniform distribution of sequences*. Courier Corpora-
725 tion, 2012.
- 726 Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- 727 Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril
728 Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, Sumith Kulal, Kyle Lacey,
729 Yam Levi, Cheng Li, Dominik Lorenz, Jonas Müller, Dustin Podell, Robin Rombach, Harry Saini,
730 Axel Sauer, and Luke Smith. Flux.1 kontext: Flow matching for in-context image generation and
731 editing in latent space, 2025. URL <https://arxiv.org/abs/2506.15742>.
- 732 Adam Daniel Laud. *Theory and application of reward shaping in reinforcement learning*. University
733 of Illinois at Urbana-Champaign, 2004.
- 734 Jongmin Lee, Wonseok Jeon, Geon-Hyeong Kim, and Kee-Eung Kim. Monte-carlo tree search in
735 continuous action spaces with value gradients. In *Proceedings of the AAAI conference on artificial*
736 *intelligence*, volume 34, pp. 4561–4568, 2020.
- 737 Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel,
738 Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human
739 feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- 740 Dexun Li, Cong Zhang, Kuicai Dong, Derrick Goh Xin Deik, Ruiming Tang, and Yong Liu.
741 Aligning crowd feedback via distributional preference reward modeling. *arXiv preprint*
742 *arXiv:2402.09764*, 2024a.
- 743 Mingxiao Li, Tingyu Qu, Ruicong Yao, Wei Sun, and Marie-Francine Moens. Alleviating ex-
744 posure bias in diffusion models through sampling with shifted time steps. *arXiv preprint*
745 *arXiv:2305.15583*, 2023.
- 746 Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xincheng Deng, Yingfang Zhang,
747 Xingchao Liu, Minbin Huang, Zedong Xiao, et al. Hunyuan-dit: A powerful multi-resolution
748 diffusion transformer with fine-grained chinese understanding. *arXiv preprint arXiv:2405.08748*,
749 2024b.

- 756 Zhanhao Liang, Yuhui Yuan, Shuyang Gu, Bohan Chen, Tiankai Hang, Ji Li, and Liang Zheng.
757 Step-aware preference optimization: Aligning preference with denoising performance at each
758 step. *arXiv preprint arXiv:2406.04314*, 2(3), 2024.
- 759
760 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
761 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European*
762 *conference on computer vision*, pp. 740–755. Springer, 2014.
- 763 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
764 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 765
766 Buhua Liu, Shitong Shao, Bao Li, Lichen Bai, Zhiqiang Xu, Haoyi Xiong, James Kwok, Sumi
767 Helal, and Zeke Xie. Alignment of diffusion models: Fundamentals, challenges, and future.
768 *arXiv preprint arXiv:2409.07253*, 2024.
- 769 Fangfu Liu, Hanyang Wang, Yimo Cai, Kaiyan Zhang, Xiaohang Zhan, and Yueqi Duan. Video-t1:
770 Test-time scaling for video generation. *arXiv preprint arXiv:2503.18942*, 2025a.
- 771
772 Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan,
773 Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *arXiv*
774 *preprint arXiv:2505.05470*, 2025b.
- 775 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and
776 transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- 777
778 Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool.
779 Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the*
780 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 11461–11471, 2022.
- 781 Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan Yang,
782 Yandong Li, Tommi Jaakkola, Xuhui Jia, et al. Inference-time scaling for diffusion models beyond
783 scaling denoising steps. *arXiv preprint arXiv:2501.09732*, 2025.
- 784 Madebyollin. madebyollin/sd-xl-vae-fp16-fix · hugging face, 2024. URL [https://](https://huggingface.co/madebyollin/sd-xl-vae-fp16-fix)
785 huggingface.co/madebyollin/sd-xl-vae-fp16-fix.
- 786
787 Oscar Mañas, Pietro Astolfi, Melissa Hall, Candace Ross, Jack Urbanek, Adina Williams, Aish-
788 warya Agrawal, Adriana Romero-Soriano, and Michal Drozdal. Improving text-to-image con-
789 sistency via automatic prompt optimization. *arXiv preprint arXiv:2403.17804*, 2024.
- 790 Jiafeng Mao, Xueting Wang, and Kiyoharu Aizawa. The lottery ticket hypothesis in denoising:
791 Towards semantic-driven initialization. In *European Conference on Computer Vision*, pp. 93–
792 109. Springer, 2024.
- 793
794 Shunqi Mao, Wei Guo, Chaoyi Zhang, Jieting Long, Ke Xie, and Weidong Cai. Ctrl-z sampling: Dif-
795 fusion sampling with controlled random zigzag explorations. *arXiv preprint arXiv:2506.20294*,
796 2025.
- 797 Farida Memon, Mukhtiar Ali Unar, and Sheeraz Memon. Image quality assessment for performance
798 evaluation of focus measure operators. *Mehran University Research Journal of Engineering &*
799 *Technology*, 34(4):379–386, 2015.
- 800
801 Wenyi Mo, Tianyu Zhang, Yalong Bai, Bing Su, Ji-Rong Wen, and Qing Yang. Dynamic prompt op-
802 timizing for text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer*
803 *Vision and Pattern Recognition*, pp. 26627–26636, 2024.
- 804 Thomas M Moerland, Joost Broekens, Aske Plaat, and Catholijn M Jonker. A0c: Alpha zero in
805 continuous action space. *arXiv preprint arXiv:1805.09613*, 2018.
- 806
807 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
808 In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.
- 809 Bernt Oksendal. *Stochastic differential equations: an introduction with applications*. Springer
Science & Business Media, 2013.

- 810 OpenCV. Github - opencv/opencv-python: Automated ci toolchain to produce precompiled opencv-
811 python, opencv-python-headless, opencv-contrib-python and opencv-contrib-python-headless
812 packages, 2024. URL <https://github.com/opencv/opencv-python>. [Online; ac-
813 cessed 2025-05-23].
- 814 Yuta Oshima, Masahiro Suzuki, Yutaka Matsuo, and Hiroki Furuta. Inference-time text-to-video
815 alignment with diffusion latent beam search. *arXiv preprint arXiv:2501.19252*, 2025.
- 816
- 817 Gaurav Parmar, Or Patashnik, Daniil Ostashev, Kuan-Chieh (Jackson) Wang, Kfir Aberman, Srimi-
818 vasa Narasimhan, and Jun-Yan Zhu. Scaling group inference for diverse and high-quality gener-
819 ation. *arXiv preprint arXiv:2508.15773*, 2025.
- 820 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
821 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
- 822
- 823 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
824 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
825 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 826 Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-
827 image diffusion models with reward backpropagation. 2023.
- 828
- 829 Zipeng Qi, Lichen Bai, Haoyi Xiong, and Zeke Xie. Not all noises are created equally: Diffusion
830 noise selection and optimization. *arXiv preprint arXiv:2407.14041*, 2024.
- 831 Kian Hong Quah and Chai Quek. Maximum reward reinforcement learning: A non-cumulative
832 reward criterion. *Expert Systems with Applications*, 31(2):351–359, 2006.
- 833
- 834 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
835 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
836 models from natural language supervision. In *International conference on machine learning*, pp.
837 8748–8763. PmLR, 2021.
- 838 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
839 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
840 *in Neural Information Processing Systems*, 36:53728–53741, 2023.
- 841 AM Raid, WM Khedr, Mohamed A El-Dosuky, and Wesam Ahmed. Jpeg image compression using
842 discrete cosine transform-a survey. *arXiv preprint arXiv:1405.6147*, 2014.
- 843
- 844 Vignav Ramesh and Morteza Mardani. Test-time scaling of diffusion models via noise trajectory
845 search. *arXiv preprint arXiv:2506.03164*, 2025.
- 846 Benjamin Rivière, John Lathrop, and Soon-Jo Chung. Monte carlo tree search with spectral expan-
847 sion for planning with dynamical systems. *Science Robotics*, 9(97):eado1010, 2024.
- 848
- 849 Christoph Schuhmann Romain Beaumont. Github - laion-ai/aesthetic-predictor: A linear esti-
850 mator on top of clip to predict the aesthetic quality of pictures. [https://github.com/](https://github.com/LAION-AI/aesthetic-predictor)
851 [LAION-AI/aesthetic-predictor](https://github.com/LAION-AI/aesthetic-predictor), 2022.
- 852
- 853 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
854 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
ence on computer vision and pattern recognition, pp. 10684–10695, 2022.
- 855
- 856 Amirmojtaba Sabour, Sanja Fidler, and Karsten Kreis. Align your steps: Optimizing sampling
857 schedules in diffusion models. *arXiv preprint arXiv:2404.14507*, 2024.
- 858
- 859 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar
860 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic
861 text-to-image diffusion models with deep language understanding. *Advances in neural informa-*
tion processing systems, 35:36479–36494, 2022.
- 862
- 863 Dvir Samuel, Rami Ben-Ari, Simon Raviv, Nir Darshan, and Gal Chechik. Generating images of
rare concepts using pre-trained diffusion models. In *Proceedings of the AAAI Conference on*
Artificial Intelligence, volume 38, pp. 4695–4703, 2024.

- 864 Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rom-
865 bach. Fast high-resolution image synthesis with latent adversarial diffusion distillation. In *SIG-*
866 *GRAPH Asia 2024 Conference Papers*, pp. 1–11, 2024a.
- 867 Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion dis-
868 tillation. In *European Conference on Computer Vision*, pp. 87–103. Springer, 2024b.
- 869 John J Shynk. *Probability, random variables, and random processes: theory and signal processing*
870 *applications*. John Wiley & Sons, 2012.
- 871 Raghav Singhal, Zachary Horvitz, Ryan Teehan, Mengye Ren, Zhou Yu, Kathleen McKeown, and
872 Rajesh Ranganath. A general framework for inference-time scaling and steering of diffusion
873 models. *arXiv preprint arXiv:2501.06848*, 2025.
- 874 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
875 *preprint arXiv:2010.02502*, 2020.
- 876 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data dis-
877 tribution. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and
878 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran
879 Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2019/file/3001ef257407d5a371a96dcd947c7d93-Paper.pdf)
880 [paper/2019/file/3001ef257407d5a371a96dcd947c7d93-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/3001ef257407d5a371a96dcd947c7d93-Paper.pdf).
- 881 Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben
882 Poole. Score-based generative modeling through stochastic differential equations, 2021. URL
883 <https://arxiv.org/abs/2011.13456>.
- 884 Zhiwei Tang, Jiangweizhi Peng, Jiasheng Tang, Mingyi Hong, Fan Wang, and Tsung-Hui Chang.
885 Inference-time alignment of diffusion models with direct noise optimization. *arXiv preprint*
886 *arXiv:2405.18881*, 2024.
- 887 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,
888 volume 47. Cambridge university press, 2018.
- 889 Grigorii Vevurko, Wendelin Böhmer, and Mathijs De Weerd. To the max: reinventing reward in
890 reinforcement learning. *arXiv preprint arXiv:2402.01361*, 2024.
- 891 Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2008.
- 892 Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Ra-
893 sul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and
894 Thomas Wolf. Diffusers: State-of-the-art diffusion models. [https://github.com/](https://github.com/huggingface/diffusers)
895 [huggingface/diffusers](https://github.com/huggingface/diffusers), 2022.
- 896 Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,
897 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using
898 direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
899 *and Pattern Recognition*, pp. 8228–8238, 2024.
- 900 Bin Xu Wang and John J Vastola. Diffusion models generate images like painters: an analytical
901 theory of outline first, details later. *arXiv preprint arXiv:2303.02490*, 2023.
- 902 Guangyi Wang, Wei Peng, Lijiang Li, Wenyu Chen, Yuren Cai, and Songzhi Su. Diffusion sampling
903 correction via approximately 10 parameters. *arXiv preprint arXiv:2411.06503*, 2024.
- 904 Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion
905 null-space model. *arXiv preprint arXiv:2212.00490*, 2022.
- 906 Yunlong Wang, Shuyuan Shen, and Brian Y Lim. Reprompt: Automatic prompt editing to refine ai-
907 generative art towards precise expressions. In *Proceedings of the 2023 CHI conference on human*
908 *factors in computing systems*, pp. 1–29, 2023.
- 909 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.
910 Human preference score v2: A solid benchmark for evaluating human preferences of text-to-
911 image synthesis. *arXiv preprint arXiv:2306.09341*, 2023a.

- 918 Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score:
919 Better aligning text-to-image models with human preference. In *Proceedings of the IEEE/CVF*
920 *International Conference on Computer Vision*, pp. 2096–2105, 2023b.
- 921 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
922 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.
923 *Advances in Neural Information Processing Systems*, 36:15903–15935, 2023a.
- 924 Yilun Xu, Mingyang Deng, Xiang Cheng, Yonglong Tian, Ziming Liu, and Tommi Jaakkola. Restart
925 sampling for improving generative processes. *Advances in Neural Information Processing Sys-*
926 *tems*, 36:76806–76838, 2023b.
- 927 Haolin Yang, Feilong Tang, Ming Hu, Yulong Li, Yexin Liu, Zelin Peng, Junjun He, Zongyuan Ge,
928 and Imran Razzak. Scalingnoise: Scaling inference-time search for generating infinite videos,
929 2025. URL <https://arxiv.org/abs/2503.16400>.
- 930 Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiabin Chen, Weihao Shen, Xiaolong Zhu, and Xiu Li.
931 Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings of*
932 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8941–8951, 2024a.
- 933 Shentao Yang, Tianqi Chen, and Mingyuan Zhou. A dense reward view on aligning text-to-image
934 diffusion with preference. *arXiv preprint arXiv:2402.08265*, 2024b.
- 935 Xingyi Yang, Daquan Zhou, Jiashi Feng, and Xinchao Wang. Diffusion probabilistic model made
936 slim. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*,
937 pp. 22552–22562, 2023.
- 938 Zilyu Ye, Zhiyang Chen, Tiancheng Li, Zemin Huang, Weijian Luo, and Guo-Jun Qi. Schedule
939 on the fly: Diffusion time prediction for faster and better image generation. *arXiv preprint*
940 *arXiv:2412.01243*, 2024.
- 941 Po-Hung Yeh, Kuang-Huei Lee, and Jun-Cheng Chen. Training-free diffusion model alignment with
942 sampling demons. *arXiv preprint arXiv:2410.05760*, 2024.
- 943 Yu Yu, Weibin Zhang, and Yun Deng. Frechet inception distance (fid) for evaluating gans. *China*
944 *University of Mining Technology Beijing Graduate School*, 3, 2021.
- 945 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable
946 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on*
947 *computer vision and pattern recognition*, pp. 586–595, 2018.
- 948 Tao Zhang, Cheng Da, Kun Ding, Kun Jin, Yan Li, Tingting Gao, Di Zhang, Shiming Xiang, and
949 Chunhong Pan. Diffusion model as a noise-aware latent reward model for step-level preference
950 optimization. *arXiv preprint arXiv:2502.01051*, 2025.
- 951 Ziyi Zhang, Li Shen, Sen Zhang, Deheng Ye, Yong Luo, Miaojing Shi, Bo Du, and Dacheng
952 Tao. Aligning few-step diffusion models with dense reward difference learning. *arXiv preprint*
953 *arXiv:2411.11727*, 2024a.
- 954 Ziyi Zhang, Sen Zhang, Yibing Zhan, Yong Luo, Yonggang Wen, and Dacheng Tao. Confronting
955 reward overoptimization for diffusion models: A perspective of inductive and primacy biases.
956 *arXiv preprint arXiv:2402.08552*, 2024b.
- 957 Zhenyu Zhou, Defang Chen, Can Wang, and Chun Chen. Fast ode-based sampling for diffusion
958 models in around 5 steps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
959 *Pattern Recognition*, pp. 7777–7786, 2024a.
- 960 Zikai Zhou, Shitong Shao, Lichen Bai, Zhiqiang Xu, Bo Han, and Zeke Xie. Golden noise for
961 diffusion models: A learning framework. *arXiv preprint arXiv:2411.09502*, 2024b.
- 962 Beier Zhu, Ruoyu Wang, Tong Zhao, Hanwang Zhang, and Chi Zhang. Distilling parallel gradients
963 for fast ode solvers of diffusion models. *arXiv preprint arXiv:2507.14797*, 2025.
- 964 Le Zhuo, Liangbing Zhao, Sayak Paul, Yue Liao, Renrui Zhang, Yi Xin, Peng Gao, Mohamed
965 Elhoseiny, and Hongsheng Li. From reflection to perfection: Scaling inference-time optimization
966 for text-to-image diffusion models via reflection tuning. *arXiv preprint arXiv:2504.16080*, 2025.
- 967
968
969
970
971

972	Contents	
973		
974	1 Introduction	1
975		
976	2 Preliminaries	2
977	2.1 DDIM	2
978	2.2 DDPO-Style Modeling	3
979	2.3 Latent Reward Shaping	3
980		
981	3 Problem Statement and Our Formulation	3
982	3.1 Problem Formulation	3
983	3.2 Augmented RODE Sampling for DDIM	4
984	3.3 MDP Modeling	6
985	3.4 Optimal Control	6
986		
987	4 Online Planning Using MCTS	6
988	4.1 Selection Policies	6
989	4.2 Beta Policy for Continuous Action Space	7
990		
991	5 Experiments	7
992	5.1 Settings and Evaluations	7
993	5.2 RODE Sampling	8
994	5.3 Comparison between Reward Shaping, MDP Modeling and Value Policies	9
995	5.4 Aligning with Aesthetics	9
996	5.5 Aligning with Semantics	9
997	5.6 Aligning with Composite Rewards	10
998	5.7 Ablations and Applications	10
999		
1000	6 Related and Conclusion	10
1001		
1002	References	10
1003		
1004	Contents	19
1005		
1006	Appendix	20
1007		
1008	A Use of Large Language Models	21
1009		
1010	B Abbreviations	21
1011		
1012	C Related	21
1013	C.1 Diffusion Model Alignment	21
1014	C.2 Inference-time Scaling for Diffusion Models	22
1015	C.3 Latent Reward Shaping	23
1016		
1017	D Supplementary to Preliminary	25
1018	D.1 Supplementary to DDPO-Style Modeling	25
1019	D.2 Supplementary to Latent Reward Shaping	26
1020	D.3 Max-Reward Control/RL	27
1021	D.4 SDE Sampling for DMs	27
1022	D.5 UCB in MCTS	28
1023		
1024	E More Discussion on Our Formulation	29
1025	E.1 Flaws of the Conventional Formulation	29
	E.2 Comparison with Existing Inference-Time Scaling Paradigms	29
	E.3 Comparison with Existing Methods	30
	E.4 Relationship to EDM Framework	32
	F More Discussion on Augmented RODE Sampling	32
	F.1 RODE Sampling for DDIM	32

1026	F.2 Variance Analysis	33
1027	F.3 Error Analysis	34
1028	F.4 Augmented RODE Sampling	37
1029	F.5 Comparison with Previous Works	38
1030		
1031	G More Discussion on Our MCTS	38
1032	G.1 Overview	38
1033	G.2 Supplementary to Mode Re-Parameterization	39
1034	G.3 Supplementary to Initialization and Online Update	39
1035	G.4 Efficiency Optimization	39
1036		
1037	H Experimental Settings, Evaluation Dimensions and Reward Functions	40
1038	H.1 Implementation	40
1039	H.2 Settings	40
1040	H.3 Metric	40
1041	H.4 Overview of Experimental Composition	43
1042		
1043	I Experiments: Supplementary to RODE Sampling	43
1044	I.1 Supplementary to Trajectory Simulation	43
1045	I.2 Supplementary to Sample Diversity	43
1046	I.3 Supplementary to Empirical Distribution	44
1047	I.4 Visualization of Sampling Trajectories	45
1048		
1049	J Experiment: Supplementary to Aligning with Aesthetics	45
1050	J.1 Settings	45
1051	J.2 Results and Analysis	46
1052		
1053	K Experiment: Supplementary to Aligning with Semantics	49
1054	K.1 Settings	49
1055	K.2 Results and Analysis	50
1056		
1057	L Experiment: Supplementary to Aligning with Composite Rewards	51
1058	L.1 Settings	51
1059	L.2 Results and Analysis	52
1060		
1061	M Experiment: Supplementary to Generalization to Other Latent Reward Policies	53
1062	M.1 Settings	53
1063	M.2 Results and Analysis	53
1064		
1065	N Experiment: Supplementary to Ablations	53
1066	N.1 Supplementary to Reward Hacking and Effects of τ	53
1067	N.2 Ablations on m and ζ	54
1068	N.3 Supplementary to Main Ablations	54
1069		
1070	O Applications	55
1071	O.1 Synergy with Community Modules	55
1072	O.2 Robustness of Image Reward Functions	56
1073	O.3 Quantitative Evaluation of Initial Noise Potential	57
1074		
1075	P Failure Cases	57
1076	P.1 Aligning with Aesthetics	57
1077	P.2 Aligning with Semantics	57
1078		
1079	Q Broader Impacts, Limitations and Future Work	57
	Q.1 Broader Impacts	57
	Q.2 Limitations	58
	Q.3 Future Work	58

Appendix

A USE OF LARGE LANGUAGE MODELS

We use ChatGPT 4.0 and Grok-Expert to aid and polish writing.

B ABBREVIATIONS

Tab. 9 presents abbreviations and their corresponding full names in this paper.

Table 9: **Abbreviations and their corresponding full names.**

Abbreviation	Full Name
DM	diffusion model
RL	reinforcement learning
MDP	Markov decision process
MCTS	Monte Carlo tree search
UCB	upper confidence bound
ODE	ordinary differential equation
SDE	stochastic differential equation
RODE	random ordinary differential equation
Aug. RODE sampling	augmented RODE sampling
NFE	number of function evaluation
MPD	mean pairwise distance
LAPV	Laplacian variance
CCR	color channel reward

C RELATED

C.1 DIFFUSION MODEL ALIGNMENT

C.1.1 HUMAN PREFERENCE QUANTIFICATION

Human preferences should be quantified first before achieving DM alignment. For example, CLIP Score (Radford et al., 2021; Hessel et al., 2021) is widely adopted to measure the semantic alignment of generated samples. However, some preferences are hard to formulate, for example, aesthetics. To overcome this challenge, researchers turn to collect preference data to train models to embed these preferences implicitly. The most representative ones are the models that predict human aesthetic preferences, such as AesScore (Romain Beaumont, 2022), PickScore (Kirstain et al., 2023), ImageReward (Xu et al., 2023a) and Human Preference Score v2 (HPS v2) (Wu et al., 2023b;a).

C.1.2 TRAINING-BASED ALIGNMENT

Training-based DM alignment techniques like Lee et al. (2023); Black et al. (2023); Fan et al. (2023) model the inference process of DMs as a multi-step MDP, and fine-tune or train DMs to align them to the reward functions (existing or trained alongside DMs), such as image aesthetics (Black et al., 2023; Fan et al., 2023; Lee et al., 2023; Xu et al., 2023a; Clark et al., 2023), semantic alignment (Fan et al., 2023; Lee et al., 2023), (in)compressibility (Black et al., 2023; Clark et al., 2023), object detection and removal (Clark et al., 2023), *et al.*

We omit to review a category of training-based DM alignment methods — direct preference optimization (DPO) (Rafailov et al., 2023; Wallace et al., 2024; Yang et al., 2024a; Li et al., 2024a; Yang et al., 2024b; Liang et al., 2024), as they lack explicit rewards, and are incompatible with our framework.

1134 C.1.3 TRAINING-FREE ALIGNMENT

1135 Training-free techniques freeze the DM’s parameters, and achieve alignment by altering its inputs
1136 or the denoising trajectory. They can be categorized into gradient-based and gradient-free methods.

1137
1138 **Gradient-based Methods.** ReNO (Eyring et al., 2024) optimizes the initial noise in a one-step
1139 inference DM using the gradient of a combination of reward objectives with respect to samples.
1140 Besides, Universal Guidance (Bansal et al., 2023) computes the gradient of the loss function with
1141 respect to latents, and applying gradient ascent to drive the intermediate latents to high-reward re-
1142 gions. Although training-free, these methods necessitate differentiable reward functions and often
1143 limited to few-step or single-step inference DMs, as gradient propagation through DMs or reward
1144 models is computationally intractable for multi-step reasoning scenarios (Xu et al., 2023a; Clark
1145 et al., 2023; Prabhudesai et al., 2023; Eyring et al., 2024).

1146
1147 **Gradient-free Methods.** To mitigate the necessity of gradient computation, researchers recently
1148 turn to gradient-free methods. For example, Z-Sampling (Bai et al., 2024) injects information via the
1149 difference in CFG (Ho & Salimans, 2022) scales between denoising and inversion to enhance both
1150 aesthetics and semantic alignment, which is not applicable to scenarios that disable CFG (e.g., un-
1151 conditional generation). Our method likewise belongs to the gradient-free family, operating without
1152 requiring differentiable reward functions or CFG.

1153 We refer readers to Liu et al. (2024) for more literature on DM alignment.

1154 C.2 INFERENCE-TIME SCALING FOR DIFFUSION MODELS

1155 Inspired by Ma et al. (2025) and Liu et al. (2024), we categorize existing DM’s inference-time
1156 scaling techniques as follows:

- 1157 1. *Inference step paradigm*: Boost performance with more inference steps, e.g., ID-
1158 DPM (Nichol & Dhariwal, 2021) and ADM-IP (Li et al., 2023);
- 1159 2. *Initial noise selection paradigm*: Use a Best-of-N strategy to pick samples with highest
1160 verifier scores from multiple initial noise candidates for better visual quality, diversity, and
1161 content consistency (for video generation), e.g., SeedSelect (Samuel et al., 2024), Scaling-
1162 Noise (Yang et al., 2025), the noise-level scaling of ReflectionFlow (Zhuo et al., 2025),
1163 Ahn et al. (2024), Qi et al. (2024), Samuel et al. (2024), Ma et al. (2025);
- 1164 3. *Initial noise construction paradigm*: Mao et al. (2024) directly creates initial noise from
1165 known winning tickets for each concept mentioned in the prompt;
- 1166 4. *Initial noise optimization paradigm*: Optimize the initial noise itself (e.g., InitNO (Guo
1167 et al., 2024), ReNO (Eyring et al., 2024), Golden Noise (Zhou et al., 2024b), DNO (Tang
1168 et al., 2024)) or its sampling distribution (e.g., InitNO (Guo et al., 2024)) via gradient-free
1169 (e.g., zero-order search (Ma et al., 2025)) or gradient-based (e.g., first-order search (Ma
1170 et al., 2025)) approaches;
- 1171 5. *Prompt optimization paradigm*: Refine the user-provided text prompts to model-preferred
1172 ones that contain richer, finer-grained descriptions (e.g., Reprompt (Wang et al., 2023),
1173 Promptist (Hao et al., 2023), OPT2I (Mañas et al., 2024), the prompt-level scaling of
1174 ReflectionFlow (Zhuo et al., 2025)), and adaptively adjust their injection timesteps and
1175 strengths (e.g., PAE (Mo et al., 2024));
- 1176 6. *Trajectory search paradigm*: Select multiple process noises as expansion directions at each
1177 timestep, and retain the top N candidates. For example, BeamDiffusion (Fernandes et al.,
1178 2025) and DBLS (Oshima et al., 2025) utilizes beam search (BS) for sampling trajectory
1179 search, where candidates are generated by sampling a set of process noise for stochastic
1180 DDIM sampling (Song et al., 2020) at each expanded node. EvoSearch (He et al., 2025)
1181 reinterprets the denoising trajectory as an evolutionary path, and mutates both the initial
1182 noise an intermediate latents for exploration. Video-T1 (Liu et al., 2025a) proposes Tree
1183 of Frame (ToF) Search to efficiently scale video generation, which adaptively expands and
1184 prunes branches in autoregressive manner. Ctrl-Z Sampling (Mao et al., 2025) conducts
1185 trajectory search with zig-zag operation (Bai et al., 2024). FK Steering (Singhal et al.,
1186 2025)

2025) scores and resamples particles according to their potentials during generation to steer them towards high-reward regions;

7. *Trajectory optimization paradigm*: Perform gradient-free information injection along the denoising trajectory (e.g., Z-Sampling (Bai et al., 2024)), or apply gradient-based guidance (e.g., Universal Guidance (Bansal et al., 2023));
8. *Correction paradigm*: The reflection-level scaling of ReflectionFlow (Zhuo et al., 2025) trains a corrector first, then uses it to progressively refine the generated images throughout the inference process.

However, we highlight the following limitations:

1. IDDPM (Nichol & Dhariwal, 2021) and ADM-IP (Li et al., 2023) found that, increasing inference steps yields diminishing FID (Yu et al., 2021) improvement, and can even degrade performance with too many steps;
2. Initial noise selection or construction methods are mainly limited to boost sample quality (e.g., FID (Yu et al., 2021) and aesthetics) and semantic alignment, which is not suitable for broader reward functions. Besides, they adopt a best-of-N procedure, i.e., independently sample multiple initial noise, and discard the sub-optimal ones along the search process. Such strategies, while acknowledging the significant influence of initial noise (Appx. I.2), neglect the contribution from the denoising trajectory, thus fail to fully exploit the generation potential in any single initial noise;
3. Initial noise optimization methods requires costly gradient calculation or approximation, limiting their deployment in multi-step inference scenarios. Besides, they also overlook trajectories;
4. Gradient-free trajectory optimization methods mainly limited to enhance the sample quality and semantic alignment (e.g., Z-Sampling (Bai et al., 2024)), which is always inapplicable to other reward functions;
5. Since initial noise, process noise (e.g., DBLS (Oshima et al., 2025)), and intermediate latents are high-dimensional vector, all methods except *inference step paradigm* adopt high-dimensional action spaces in MDP contexts, leading to low sampling efficiency and high variance.

To overcome the above challenges, we first model the DM’s inference process as a MDP with dense rewards, then fix the initial noise and perform efficient trajectory search with MCTS (Coulom, 2006). Our method belongs to the *trajectory search paradigm*. Besides, existing inference-time scaling approaches that adopt high-dimensional action spaces (e.g., initial noise, process noise, intermediate latents) suffer from low parameter efficiency (Sec. 5.4, Appx. J). Recent works (Zhou et al., 2024a; Chen et al., 2023a; 2025a; Wang et al., 2024; Zhu et al., 2025) designate novel distillation techniques with minimal learnable parameters, which inspire us to adopt an extremely low-dimensional action space — 1-D action space — for parameter-efficient scaling.

Besides, Parmar et al. (2025) marks a first step towards scaling a group of samples for better diversity and quality.

Comparison. Tab. 10 compares our method against some existing training-based and training-free methods (especially trajectory-based paradigms). Note that our method requires neither differentiable reward functions nor gradient propagation, which is applicable to a wide range of reward functions.

C.3 LATENT REWARD SHAPING

Most existing alignment techniques adopt the DDPO-style modeling (Black et al., 2023) (Sec. 2.2), which suffers from the sparse reward issue in reinforcement learning (RL), i.e., the reward signals emerge only at the end of the process, which causes unstable training, slow convergence, and inefficient data use (Hare, 2019). On the one hand, adopting sparse rewards lead the training or guidance process to focus too much on the rewards of the final samples, neglecting the contribution from the denoising trajectory. On the other hand, in the practices of existing RL-based DM alignment, methods using dense rewards (Zhang et al., 2024b;a; Ye et al., 2024; Zhang et al., 2025) often outperform

Table 10: **Comparison between popular DM alignment methods and ours.** “*Train.*” denotes whether training or fine-tuning is required; “*Dense.*” denotes whether the reward signals are dense; “*Diff.*” denotes whether necessitating differentiable reward functions; “*Gradient.*” denotes whether gradient propagation through the DMs or the reward models is needed beyond training. Methods that are more user-friendly (less restrictive, more flexible, more computationally tractable) are indicated in green. Particularly, ✓ denotes the need to train both the DM and an additional model; ✓ denotes the need to train an additional model rather than the DM itself; ✓ indicates dense rewards due to single-step inference.

Methods	Train.	Dense.	Diff.	Gradient.
DDPO (Black et al., 2023)	✓	✗	✗	✗
DPOK (Fan et al., 2023)	✓	✗	✗	✗
Lee et al. (Lee et al., 2023)	✓	✗	✗	✗
ReFL (Xu et al., 2023a)	✓	✓	✓	1 step
DRaFT-K (Clark et al., 2023)	✓	✓	✓	fixed K steps
AlignProp (Prabhudesai et al., 2023)	✓	✓	✓	random steps
TDPO-R (Zhang et al., 2024b)	✓	✓	✗	✗
SDPO (Zhang et al., 2024a)	✓	✓	✗	✗
Schedule On the Fly (Ye et al., 2024)	✓	✓	✗	✗
LRM (Zhang et al., 2025)	✓	✓	✗	✗
ReNO (Eyring et al., 2024)	✗	✓	✓	1 step
DNO (Tang et al., 2024)	✗	✗	✓	n step
Universal Guidance (Bansal et al., 2023)	✗	✓	✓	multiple steps
DBLS (Oshima et al., 2025)	✗	✓	✗	✗
EvoSearch (He et al., 2025)	✗	✓	✗	✗
Video-T1 (Liu et al., 2025a)	✗	✓	✗	✗
Ctrl-Z Sampling (Mao et al., 2025)	✗	✓	✗	✗
FK Steering (Singhal et al., 2025)	✗	✓	✗	✗
Ours	✗	✓	✗	✗

those with sparse rewards (Black et al., 2023; Fan et al., 2023; Lee et al., 2023; Xu et al., 2023a; Clark et al., 2023; Prabhudesai et al., 2023). We attribute it to the fact that, under the same computational resources, methods with dense rewards receive more supervision signals in the process.

Recently, researchers have tried designing intermediate rewards for the multi-step DM inference process to obtain dense supervision signals (Zhang et al., 2024b;a; Dhariwal & Nichol, 2021), which is called *reward shaping* (Laud, 2004; Grzes, 2017) in the RL context. The challenge is that, most popular reward models for image synthesis (e.g., FID (Yu et al., 2021), IS (Barratt & Sharma, 2018), LPIPS (Zhang et al., 2018), aesthetics prediction models (Kirstain et al., 2023; Xu et al., 2023a; Romain Beaumont, 2022; Wu et al., 2023b;a)) are trained on clean image domains, which can not be directly transferred to the noisy latent domain to provide reliable intermediate rewards. Although ImageReward (Xu et al., 2023a) found that, rewards calculated from noisy images after 30-th denoising step in a 40-step inference are reliable enough, obtaining rewards for every denoising step, especially in the early stages, is still tough (Xu et al., 2023a; Clark et al., 2023; Prabhudesai et al., 2023; Eyring et al., 2024).

There are training-based and training-free methods that focus on deriving intermediate rewards for the noisy latent domain from reward models trained on clean image domains (we call them *clean reward models*), which we call *latent reward shaping*.

Training-based Methods. An immediate thoughts is to train an extra counterpart of the clean reward model in the latent domain to accept noisy inputs, such as classifier guidance (Dhariwal & Nichol, 2021) and LRM (Zhang et al., 2025). Recently, discount-based methods have been developed for latent reward acquisition. For example, TDPO-R (Zhang et al., 2024b) trains an extra model to discount the final rewards into intermediate rewards. Besides, vision language models (VLMs) have also be adopted as reward critics for intermediate rewards (Liu et al., 2025a). However, these

1296 methods demand lots of training, may not generalize across datasets, and require retraining when
 1297 reward function changes.
 1298

1299 **Training-free Methods.** To overcome the limitations of training-based approaches, recent works
 1300 design training-free paradigms to efficiently obtain intermediate rewards. The key idea of training-
 1301 free latent reward shaping is to use the final reward as a *proxy* for intermediate rewards. For example,
 1302 Universal Guidance (Bansal et al., 2023) computes guidance using the *pseudo-final samples* derived
 1303 from the posterior mean of DDIM (Song et al., 2020). DBLS (Oshima et al., 2025) uniformly inter-
 1304 polates the remaining denoising steps into 2 or 3 DDIM steps for more accurate estimation. Schedule
 1305 On the Fly (Ye et al., 2024) discounts the final rewards backward to obtain intermediate rewards.
 1306 SDPO (Zhang et al., 2024a) obtains dense rewards via the cosine similarity between the intermedi-
 1307 ate latents and those at the initial, final, or other anchor steps. We present a general formulation for
 1308 training-free latent reward shaping in Sec. 2.3 in the main paper and Appx. D.2.
 1309

1310 D SUPPLEMENTARY TO PRELIMINARY

1311 D.1 SUPPLEMENTARY TO DDPO-STYLE MODELING

1312 Given a timestep schedule $\mathcal{T} = \{t_0, t_1, \dots, t_n\}$, where $T = t_0 > t_1 > \dots > t_n = 0$, in
 1313 which T represents the number of training steps of the DM, and n denotes the number of inference
 1314 steps. Meanwhile, given a DM p_θ parameterized by θ , and a clean reward model $\phi(\cdot)$. Condition
 1315 information (such as text prompts) is omitted for simplicity. The inference process of the DM is
 1316 modeled as an MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma, K)$, where:
 1317

1318 (a) The state space $\mathcal{S} \subseteq \mathbb{R}^{c \times h \times w} \times \mathcal{T}$ consists of all timesteps t_i ($0 \leq i \leq n$) and their corresponding
 1319 latent variables z_{t_i} with varying noise levels. Here, c, h, w are the number of channels, height, and
 1320 width of the latent variables, respectively. The state at timestep t_i is defined as
 1321

$$1322 s_{t_i} \triangleq (z_{t_i}, t_i); \quad (18)$$

1323 (b) The action space $\mathcal{A} \subseteq \mathbb{R}^{c \times h \times w}$ comprises all latent variables z_{t_i} ($0 \leq i \leq n$) across noise levels.
 1324 The action at timestep t_i is defined as the latent variable corresponding to the noise level of the next
 1325 timestep t_{i+1} . The policy and action are given by
 1326

$$1327 \pi(a_{t_i}, s_{t_i}) \triangleq p_\theta(z_{t_{i+1}} | z_{t_i}, t_i), \quad a_{t_i} \triangleq z_{t_{i+1}}; \quad (19)$$

1328 (c) The transition dynamics P is determined by p_θ . Formally,
 1329

$$1330 P(s_{t_{i+1}} | s_{t_i}, a_{t_i}) \triangleq (\delta_{z_{t_{i+1}}}, \delta_{t_{i+1}}), \quad (20)$$

1331 in which δ_y is the Dirac delta distribution with non-zero density only at y . The specific form when
 1332 using DDIM (Song et al., 2020) writes Eq. 1 in Sec. 2.1;
 1333

1334 (d) The reward function is defined as
 1335

$$1336 R(s_{t_i}, a_{t_i}, s_{t_{i+1}}) \triangleq \begin{cases} \mathbb{E}_{z_{t_n} \sim p_\theta(z_{t_{i+1}})}[\phi(z_{t_n})] & \text{if } i + 1 = n \\ 0 & \text{otherwise} \end{cases}; \quad (21)$$

1337 (e) The discount factor γ is within the interval $(0, 1]$. Since rewards are only obtainable at the final
 1338 timestep (Eq. 21), γ becomes useless in DDPO-style modeling. We still explicitly include it here to
 1339 maintain consistency with subsequent representations;
 1340

1341 (f) The time horizon K is defined as the number of inference steps, *i.e.*,

$$1342 K \triangleq n, \quad (22)$$

1343 which is incorporated into the MDP tuple to emphasize that it’s an episodic MDP (Grzes, 2017).
 1344

D.2 SUPPLEMENTARY TO LATENT REWARD SHAPING

D.2.1 VARIANTS OF \mathcal{F} FOR LATENT REWARD SHAPING

This section presents several variations of the estimation policy \mathcal{F} used for deriving the *pseudo-final latent* $\hat{z}_{t_n} = \mathcal{F}(z_{t_i})$ from an intermediate latent z_{t_i} ($0 \leq i < n$).

(a) *immediate-ddim*: The posterior mean in the DDIM (Song et al., 2020) step

$$\mathbb{E}[z_{t_n} | z_{t_i}] = \frac{z_{t_i} - \sqrt{1 - \bar{\alpha}_{t_i}} \epsilon_{\theta}^{(t_i)}(z_{t_i})}{\sqrt{\bar{\alpha}_{t_i}}} \quad (23)$$

can be used to derive pseudo-final latents, *i.e.*,

$$\mathcal{F}(z_{t_i}) = \mathbb{E}[z_{t_n} | z_{t_i}]; \quad (24)$$

(b) *immediate-score*: Similar to (a) that directly predicts \hat{z}_{t_n} from z_{t_i} , not assuming z_{t_i} is noise-free, Chung et al. (2022) propose a more accurate posterior mean estimation based on the score function (Song & Ermon, 2019), which writes

$$\mathbb{E}[z_{t_n} | z_{t_i}] = \frac{z_{t_i} + (1 - \bar{\alpha}_{t_i}) \nabla_{z_{t_i}} \log p_{t_i}(z_{t_i})}{\sqrt{\bar{\alpha}_{t_i}}} \approx \frac{z_{t_i} + (1 - \bar{\alpha}_{t_i}) s_{\theta}(z_{t_i})}{\sqrt{\bar{\alpha}_{t_i}}}, \quad (25)$$

where $s_{\theta}(z_{t_i}) = -\frac{\epsilon_{\theta}(z_{t_i})}{\sqrt{1 - \bar{\alpha}_{t_i}}}$ is the score function at time step t_i . The form of \mathcal{F} remains as in Eq. 24;

(c) *look-ahead*: Unlike immediately estimating \hat{z}_{t_n} from z_{t_i} , DBLS (Oshima et al., 2025) introduces a look-ahead (LA) estimator. Specifically, it equally interpolates the remaining timesteps from t_i to $t_n = 0$ to 2 or 3 steps, and denoises with 2 or 3 deterministic DDIM steps to obtain the pseudo-final latent. We name them *LA-2* and *LA-3*, respectively.

Formally, denote the DDIM step from timestep t_i to t_{i+1} (Eq. 1 in Sec. 2.1) as $z_{t_{i+1}} = \mathcal{G}(z_{t_i}, \epsilon_{t_i}, \{t_i, t_{i+1}\})$, then the LA-2’s derivation of pseudo-final latents can be expressed as

$$\mathcal{F}(z_{t_i}) = \mathcal{G}(z_{t_i}, \epsilon_{\theta}^{(t_i)}(z_{t_i}), \epsilon_{t_i}, \{t_i, t_j, t_n\}), \quad (26)$$

where $t_j = \left\lfloor \frac{t_i + t_n}{2} \right\rfloor = \left\lfloor \frac{t_i}{2} \right\rfloor$, while the LA-3’s writes

$$\mathcal{F}(z_{t_i}) = \mathcal{G}(z_{t_i}, \epsilon_{\theta}^{(t_i)}(z_{t_i}), \epsilon_{t_i}, \{t_i, t_{j_0}, t_{j_1}, t_n\}), \quad (27)$$

where $t_{j_0} = \left\lfloor \frac{2t_i}{3} \right\rfloor$, $t_{j_1} = \left\lfloor \frac{t_i}{3} \right\rfloor$;

(d) *sequential*: Starting from z_{t_i} , apply deterministic DDIM steps sequentially to obtain \hat{z}_{t_n} , *i.e.*,

$$\mathcal{F}(z_{t_i}) = \mathcal{G}(z_{t_i}, \epsilon_{\theta}^{(t_i)}(z_{t_i}), \epsilon_{t_i}, \{t_i, \dots, t_n\}). \quad (28)$$

When using *latent rewards* in multi-step inference, the number of intermediate reward queries far exceeds the number of final ones (Zhang et al., 2024a), regardless of whether training-based or training-free. Besides, the intermediate rewards derived from latent rewards suffer from inaccuracy in early denoising (Zhang et al., 2024a).

Clarification on novelty. We do **not** propose a new intermediate-reward estimator. Instead, we introduce a unified latent-reward-shaping formulation that encompasses a broad class of existing ODE-based value/intermediate-reward estimation methods. The mapping \mathcal{F} is modular and replaceable, allowing practitioners to choose an appropriate instantiation based on computational or task requirements.

D.2.2 LATENT REWARD SHAPING IN DIFFERENT STYLE

Differently, SDPO (Zhang et al., 2024a) obtains dense rewards via the cosine similarity between intermediate latents and those at the initial, final, or anchor steps. This policy can not be represented as the \mathcal{F} -style latent reward shaping in Sec. D.2.1.

We do **not** adopt this policy, because intermediate rewards should be estimated in a trajectory agnostic manner during the search process in our scenario.

1404 D.2.3 DISCOUNTED vs. UNDISCOUNTED
1405

1406 When using dense rewards, the selection of γ assigns different preferences to the optimization of
1407 cumulative returns: (i) $\gamma = 1$ treats each denoising step equally (Zhang et al., 2024b); (ii) $\gamma < 1$
1408 emphasizes the contributions from earlier (Zhang et al., 2024a) or later (Ye et al., 2024) denoising
1409 steps, depending on the discounting direction.

1410
1411 D.3 MAX-REWARD CONTROL/RL

1412 Unlike typical control/RL problems that aim to maximize the (discounted) cumulative returns,
1413 the *max-reward control/RL problem* focuses on maximizing the (discounted) rewards encountered
1414 throughout the process (**not** just the reward at the end of the episode).

1415
1416 **Deterministic Scenarios.** Quah & Quek (2006); Gottipati et al. (2020) introduces a Q-function
1417 for *deterministic* scenarios, which recursively optimizes for the expectation of maximum of reward
1418 at the current timestep and future expectations, expressed as

1419
1420
$$Q^\pi(s_t, a_t) = \mathbb{E}_{a_{t+1}}^{s_{t+1}} \left[\max \left\{ r_{t+1}, \gamma \mathbb{E}_{a_{t+2}}^{s_{t+2}} [\max\{r_{t+2}, r_{t+3}, \dots\}] \right\} \right], \quad (29)$$

1421
1422 in which $r_{t+1} \triangleq R(s_t, a_t, s_{t+1})$. This Q-function is proven to satisfy the Bellman-like equation

1423
1424
$$Q^\pi(s_t, a_t) = \mathbb{E}_{a_{t+1}}^{s_{t+1}} [\max\{r_{t+1}, \gamma Q^\pi(s_{t+1}, a_{t+1})\}]. \quad (30)$$

1425
1426 Although we do **not** employ stochastic control/RL in this study, we still present a brief introduction
1427 to highlight the simplicity of deterministic counterpart.

1428
1429 **Stochastic Scenarios.** Veviuurko et al. (2024) further extends this concept to *stochastic* environ-
1430 ments. Firstly, the *max-reward return*

1431
1432
$$G_t = \max\{r_{t+1}, \gamma r_{t+2}, \gamma^2 r_{t+3}\} \quad (31)$$

1433 in Quah & Quek (2006); Gottipati et al. (2020) is expanded to

1434
1435
$$\mathbb{E}_\pi[G_t] = \mathbb{E}_\pi[r_{t+1} \vee \gamma G_{t+1}], \quad (32)$$

1436 where and $a \vee b \triangleq \max\{a, b\}$ does not commute with \mathbb{E} .

1437 Then, with the help of an auxiliary real variable $y \in \mathbb{R}$ that propagates information about the past
1438 rewards, the *max-reward value functions* are defined as

1439
1440
$$V^\pi(s, y) = \mathbb{E}_\pi[y \vee G_t \mid s_t = s], \quad (33)$$

1441
1442
$$Q^\pi(s, a, y) = \mathbb{E}_\pi[y \vee G_t \mid s_t = s, a_t = a], \quad (34)$$

1443 which are subject to the Bellman-like equations

1444
1445
$$V^\pi(s, y) = \gamma \mathbb{E}_{a_{t+1}}^{s_{t+1}} [y' \vee V^\pi(s_{t+1}, y') \mid s_t = s], \quad (35)$$

1446
1447
$$Q^\pi(s, a, y) = \gamma \mathbb{E}_{a_{t+1}}^{s_{t+1}} [y' \vee Q^\pi(s_{t+1}, a_{t+1}, y') \mid s_t = s, a_t = a], \quad (36)$$

1448 where $y' \triangleq \frac{R(s, a, s_{t+1}) \vee y}{\gamma}$.

1449
1450
1451 D.4 SDE SAMPLING FOR DMS

1452
1453 D.4.1 SDE SAMPLING FOR DDIM

1454 The denoising process in DMs is governed by a reverse-time SDE (Song et al., 2021)

1455
1456
$$dz = [f(z, t) - g(t)^2 \nabla_z \log p_t(z)] dt + g(t) d\bar{w}, \quad (37)$$

1457 where \bar{w} is an inverse Wiener process.

For standard stochastic DDIM (Song et al., 2020) (Eq. 23),

$$f(z, t) = \frac{d}{dt} (\log \sqrt{\bar{\alpha}_t}) z, \quad g(t) = \sigma_t, \quad (38)$$

where

$$\sigma_{t_i} = \sigma_{t_i}(\eta_{t_i}) \triangleq \eta_{t_i} \omega_{t_i} \triangleq \eta_{t_i} \sqrt{\frac{1 - \bar{\alpha}_{t_{i-1}}}{1 - \bar{\alpha}_{t_i}}} \sqrt{1 - \frac{\bar{\alpha}_{t_i}}{\bar{\alpha}_{t_{i-1}}}} \quad (39)$$

is the variance of the process noise in Eq. 37, and $\eta_{t_i} \in [0, 1]$ is sampled from a uniform distribution on $[0, 1]$ or a unimodal Beta distribution to introduce randomness, independent at each timestep.

In standard SDE sampling for DDIM, the process noise ϵ_{t_i} is re-sampled at each timestep, corresponding to the $d\bar{w}$ term in Eq. 37, which provides isotropic and uncorrelated perturbation.

D.4.2 VARIANCE ANALYSIS

Consider a discrete analog of the continuous SDE sampling. Using the Euler-Maruyama method (Oksendal, 2013), the discrete-time approximation of Eq. 37 is

$$z_{t_{i+1}} = z_{t_i} + [f(z_{t_i}, t_i) - g(t_i)^2 \nabla_z p_{t_i}(z_{t_i})] \Delta t + g(t_i) \sqrt{\Delta t} \xi_{t_i}, \quad (40)$$

where $\xi_{t_i} \sim \mathcal{N}(0, I)$, and $g(t_i) \sqrt{\Delta t} \xi_{t_i}$ approximates the $d\bar{w}_t$ over the interval $[t_i, t_{i+1}]$.

The conditional expectation

$$\mathbb{E}[z_{t_{i+1}} | z_{t_i}] = z_{t_i} + [f(z_{t_i}, t_i) - g(t_i)^2 \nabla_z p_{t_i}(z_{t_i})] \Delta t. \quad (41)$$

The conditional variance

$$\begin{aligned} \text{Var}_{\text{SDE}} &= \text{Var}(z_{t_{i+1}} | z_{t_i}) \\ &= \mathbb{E} \left[\left(g(t_i) \sqrt{\Delta t} \xi_{t_i} \right) \left(g(t_i) \sqrt{\Delta t} \xi_{t_i} \right)^\top \middle| z_{t_i} \right] \\ &= g(t_i)^2 \Delta t I = \omega_{t_i}^2 \eta^2 \Delta t I, \end{aligned} \quad (42)$$

where $\eta_{t_i} = \eta$ ($0 \leq i \leq n - 1$), and $\mathbb{E}[\xi_{t_i} \xi_{t_i}^\top] = I$.

This variance is a full-rank matrix, which reflects the isotropic perturbation in different directions, and ensures diverse sample generation. Besides, it scales linearly with Δt in discrete terms.

D.4.3 ERROR ANALYSIS

Xu et al. (2023b) points out that, the *discretization error* and the *approximation error* are the primary errors involved in the SDE sampling for DDIM. It demonstrates that, the diffusive nature of the stochasticity in standard SDE sampling contracts accumulated errors due to its isotropic noise, i.e., it helps to “reset” or “forget” their cumulative impact. This contraction effect allows SDE samplers to achieve better sample quality in the large NFE regime by mitigating error accumulation across steps.

D.5 UCB IN MCTS

UCB Definition. MCTS (Coulom, 2006; Browne et al., 2012; Rivière et al., 2024) is a heuristic search algorithm for optimizing decision-making processes, where the upper confidence bound (UCB) during the selection phase for a MDP can be written as

$$\text{UCB}(s, a) = \text{exploitation}(s, a) + \lambda_{\text{exploration}} \cdot \text{exploration}(s, a), \quad (43)$$

where s is the state, a is an action used for expansion, and $\lambda_{\text{exploration}} > 0$ is a hyper-parameter used to balance exploration and exploitation.

In conventional implementations, the exploitation and exploration term are respectively given as

$$\text{exploit}(s, a) = Q(s, a), \quad \text{explore}(s, a) = \sqrt{\frac{\ln N(s)}{N(s, a)}}, \quad (44)$$

where $N(s)$ is the number of times state s is visited, and $N(s, a)$ is the number of times action a is chosen for expansion.

Average-Value Policies. In cumulative-reward RL, the $V(s)$ is typically estimated with the average return from all simulations passing through s . Formally,

$$V(s) = \left(\sum_{i=0}^{N(s)-1} R_i \right) / N(s). \quad (45)$$

Let R be the cumulative or maximum return from the simulation phase, then $V(s)$ is updated as

$$V(s) \leftarrow V(s) + \frac{R - V(s)}{N(s)}, \quad (46)$$

during the backpropagation phase. We refer to this value policy as *average-value policy*.

Dealing with Continuous Action Spaces. Standard MCTS only supports discrete action spaces, which are unavailable in our scenario, since two points in continuous spaces are different with probability 1 (Shynk, 2012). It causes $N(s, a) = 1$ almost surely and the exploration term becomes $\sqrt{\ln N(s)}$ consistently, which renders the UCB formula ineffective for balancing exploitation and exploration among actions, as it relies heavily on visit counts to under-explored options.

Recent works have employed progressive widening (Coulom, 2007; Couëtoux et al., 2011) to tackle continuous action spaces (Moerland et al., 2018; Bianchi et al., 2023; Lee et al., 2020) (e.g., DTS (Jain et al., 2025)). However, they only sample candidate actions for selection, which limits the exploration across the entire action space, and are incompatible with our stochastic Beta policy. Instead, we propose a value-based UCB to overcome this challenge (Sec. 4.1), together with Beta policies for exploration (Sec. 4.2).

E MORE DISCUSSION ON OUR FORMULATION

E.1 FLAWS OF THE CONVENTIONAL FORMULATION

We argue that, the *max-final reward formulation* is sub-optimal when adopting dense rewards, regardless of whether the final rewards are optimized directly or indirectly:

- *Indirectly*: Use latent rewards as dense rewards to train or fine-tune DMs by optimizing the cumulative rewards. It is based on the insight that, greater cumulative rewards often lead to greater final rewards. However, latent rewards usually have similar scales to the final rewards, potentially overwhelming the final reward, especially when n is large;
- *Directly*: Using latent rewards as a heuristic to guide the optimization process of the final rewards like FK Steering (Singhal et al., 2025) could potentially mitigate the above misleading. However, the *admissibility* of latent rewards as a valid heuristic is not guaranteed, since pseudo-final samples may achieve higher rewards than the terminal samples, i.e., over-estimation exists. This can mislead the optimization process into over-exploiting paths that appear promising but are actually suboptimal for final reward maximization.

E.2 COMPARISON WITH EXISTING INFERENCE-TIME SCALING PARADIGMS

Existing inference-time scaling paradigms for DMs that adopt high-dimensional action space (Appx. C.2) have several limitations:

1. Adopting the initial noise z_{t_0} or process noise ϵ_{t_i} as actions makes a high-dimensional action space $\mathcal{A} \subseteq \mathbb{R}^{c \times h \times w}$, which is hard to enumerate efficiently. Here, c, h, w are the number of channels, height and width of the latent variables, respectively;
2. High-dimensional spaces are sparse, making effective coverage difficult with limited samples. Quantitatively, sampling M points from \mathbb{R}^d faces the sphere-covering problem (Ver-shynin, 2018), that is, covering a ball of radius r requires $M \propto \left(\frac{r}{\delta}\right)^d$ samples, where δ is the covering precision. Even for small r and δ , the required M is impractically large due to the high dimension $d = c \times h \times w$;

- 1566 3. The deviation between the ϵ_{t_i} s resampled at each timestep is relatively large, leading to
 1567 high variance due to variant perturbation direction at the same timesteps in different runs,
 1568 which results in slow convergence;
 1569
 1570 4. The manifold hypothesis (Fefferman et al., 2016) posits that, high-dimensional data lies
 1571 near a low-dimensional manifold. Thus many parameters of the high-dimensional actions
 1572 may be redundant.

1573 In contrast, our RODE-based inference-time scaling paradigm (Sec. 3.4 in the main paper) offers the
 1574 following advantages:

- 1575 1. Our 1-D actions are easy to sample with 1-D probability density functions;
 1576
 1577 2. It becomes convenient to visualize the high-dimensional denoising trajectories using polar
 1578 plots (Sec. 5.2, Appx. I.4) by parameterizing on low-dimensional action sequences;
 1579
 1580 3. Our low-dimensional compact action space $\mathcal{A} = [0, 1]$ can be efficiently covered with
 1581 limited samples. Specifically, sampling M points $\{x^{(i)}\}_{i=0}^{M-1}$ in $[0, 1]$ ensures that, for any
 1582 $x^* \in [0, 1]$, there exists $x^{(j)}$ s.t. $|x^* - x^{(j)}| \leq \frac{1}{2M}$ (Kuipers & Niederreiter, 2012);
 1583
 1584 4. The standard deviation $\sigma_{t_i} = \eta_{t_i} \cdot \omega_{t_i}$ in DDIM (Eq. 3) is a linear function of η_{t_i} , where
 1585 $\omega_{t_i} = \sqrt{\frac{1 - \bar{\alpha}_{t_{i-1}}}{1 - \bar{\alpha}_{t_i}}} \sqrt{1 - \frac{\bar{\alpha}_{t_i}}{\bar{\alpha}_{t_{i-1}}}}$ is a time-variant constant, making changes in η_{t_i} smoothly
 1586 adjust σ_{t_i} . Fixing the process noise $\{\epsilon_{t_i}\}_{i=0}^{n-1}$, the perturbation direction at each timestep
 1587 is fixed, with its strength controlled by η_{t_i} . In other words, our modeling decouples noise
 1588 direction (ϵ_{t_i}) and magnitude (η_{t_i}), enabling controllable trajectory modulation and smooth
 1589 navigation of the continuous action space;
 1590
 1591 5. Our low-dimensional actions make it easy to optimize the likelihood of sampling a certain
 1592 action and the actions nearby (Sec. 4.2).

1593 E.3 COMPARISON WITH EXISTING METHODS

1594
 1595 **Comparison with an Existing Method that Frames the Data Generation as a Control Problem.** Previously, Fan & Lee (2023) frames the data generation as a control problem, and fine-tunes DDPM (Ho et al., 2020) samplers to find a shorter path to the data distribution, instead of strictly follow the DM’s backward process. Our method differs in these aspects:

- 1600 1. We only augment each timestep of the original backward process with a shortcut to the
 1601 terminal step, rather than identifying a completely novel denoising path;
 1602
 1603 2. We align pretrained DMs in a training-free manner by optimizing the action sequence
 1604 $\{\eta_{t_i}\}_{i=0}^{n-1}$ online, rather than tuning the DM parameters.

1605 Besides, treating pseudo-final samples as valid samples in our max-reward formulation, can be
 1606 viewed as a novel implementation for deviating from the DM’s pre-defined trajectory for better
 1607 performance.

1608
 1609 **Comparison with an Existing Method that Leverage MCTS for Trajectory Search.** Concurrently, Ramesh & Mardani (2025) models the denoising process as a MDP with a terminal reward, and perform trajectory search using MCTS with contextual bandits. Our method differs in these aspects:

- 1610 1. We scale low-dimensional η_{t_i} actions (Sec. 3.2) rather than the high-dimensional ϵ_{t_i} actions
 1611 employed by Ramesh & Mardani (2025);
 1612
 1613 2. We adopt dense rewards throughout the entire denoising process, rather than the solely final
 1614 rewards employed by Ramesh & Mardani (2025);
 1615
 1616 3. We adopt dense rewards for the MDP, and model DM alignment as an max-encountered-
 1617 reward optimal control problem (Sec. 3.3, Sec. 3.4) rather than the max-final reward prob-
 1618 lem in Ramesh & Mardani (2025);
 1619

- 1620 4. [Ramesh & Mardani \(2025\)](#) discretizes the continuous action space by pre-fixing a candidate
 1621 set of process noise for every timestep, whereas ours directly explores the continuous action
 1622 space via Beta polices.
 1623

1624 **Comparison with Diffusion Tree Sampling.** Also concurrently, Diffusion Tree Sampling
 1625 (DTS) ([Jain et al., 2025](#)) introduce a MCTS-based approach that samples from the reward-aligned
 1626 target density by propogating terminal rewards back through the diffusion chain, and progressively
 1627 refining value estimation. Our method differs in these aspects:

- 1628 1. We select candidate tree nodes for expansion based on our proposed value-based UCB
 1629 (Sec. 4.1), rather than the Boltzmann distribution employed by DTS;
 1630 2. We adopt a low-dimensional action space (Sec. 3.2), and expand nodes with actions sam-
 1631 pled from onlne-updated Beta distributions (Sec. 4.2), instead of directly sampling from
 1632 high-dimensional distributions as in DTS. Besides, we do **not** employ progressive widen-
 1633 ing ([Coulom, 2007](#); [Couëtoux et al., 2011](#)) to deal with continuous action spaces as done
 1634 in DTS;
 1635 3. Beyond terminal rewards, we compute intermediate rewards at every tree node (Sec. 3.3),
 1636 and introduce max-reward modeling to fully exploit these intermediate supervision signals
 1637 (Sec. 3.1);
 1638 4. We update nodes’s values with our proposed max-value policy (Sec. 4.1), which differs
 1639 from the soft value estimation adopted in DTS.
 1640

1641 **Comparison with Demons.** Demons ([Yeh et al., 2024](#)) also seems to involve scaling scalars. Our
 1642 method differs in these aspects:
 1643

- 1644 1. **Difference in the role of scalar variables.** (I) After randomly sampling a set of candidate
 1645 noise vectors $\{z^{(k)}\}_{k=1}^K$, Demons forms a new noise vector z^* by weighting them with
 1646 a set of scalars $\{b_k\}_{k=1}^K$ (Sec. 4.2 in [2]). Importantly, each b_k is computed adaptively
 1647 according to the relative advantage of $z^{(k)}$ among $\{z^{(k)}\}_{k=1}^K$ rather than from an optimiza-
 1648 tion process; (II) Differently, we pre-sample and fix a set of process-noise directions, and
 1649 the MCTS optimization explicitly searches for the scalar sequence $\{\eta_{t_i}\}_{i=0}^{n-1}$. These scalars
 1650 in our formulation are **not** adaptively computed, instead, they are variables that *being opti-*
 1651 *mized*;
 1652 2. **Clarification regarding RODE sampling.** (I) The candidate noise set of Demons
 1653 $\{z^{(k)}\}_{k=1}^K$ is resampled at every expansion step. From a continuous-time perspective, this
 1654 repeated re-sampling induces a stochastic term — a Wiener increment dw — when writ-
 1655 ing its differential-equation formulation. Consequently, the continuous-time interpretation
 1656 of Demons corresponds to an SDE sampling process; (II) Differently, we pre-sample and
 1657 fix a set of process-noise realizations and reuse them during search instead of resampling
 1658 (Appx. F.3.4). This replaces the stochastic term $g(t)dw$ with an anisotropic drift term
 1659 $\omega(t)\eta(t)h(t)dt$ (Eq. 8). With process noise fixed, both $\omega(t)$ and $h(t)$ are deterministic,
 1660 and the only stochastic component is the scalar-valued process $\eta(t)$. Consequently, the re-
 1661 sulting dynamics contain no Wiener-process term, making the formulation a random ODE
 1662 (RODE) rather than an SDE. To the best of our knowledge, prior works have not formalized
 1663 sampling via pre-fixed process-noise realizations as a RODE for inference-time alignment.
 1664

1665 Besides, our modeling decouples noise direction (ϵ_{t_i}) and magnitude (η_{t_i}), enabling controllable
 1666 trajectory modulation and smooth navigation of the continuous action space, while Demons ([Yeh](#)
 1667 [et al., 2024](#)) does **not**.
 1668

1669 **Comparison with methods that handle non-differentiable rewards.** DNO ([Tang et al., 2024](#))
 1670 does explicitly involve gradients of the reward functions, but it adopts gradient estimation to relax
 1671 the constraint of differentiability. In contrast, our method does **not** necessitate differentiable reward
 1672 functions because it is entirely simulation-based and *does not incorporate any reward gradients*.
 1673

An intuitive way to understand this distinction: Gradient-based navigation tends to be more bene-
 1674 ficial or even necessary to efficiently explore high-dimensional action spaces. However, our action

space is extremely low-dimensional, allowing MCTS to cover it effectively with a limited search budget (Appx. E.2). As a result, gradient computation or estimation is unnecessary for improving navigation efficiency. Moreover, involving gradient estimation in such a low-dimensional space may waste NFE dynamics and NFE reward evaluations, potentially leading to suboptimal performance.

E.4 RELATIONSHIP TO EDM FRAMEWORK

Solving for the optimal control sequence $\{\eta_{t_i}\}_{i=0}^{n-1}$ is similar to indirectly performing an adaptive search for a time-variant S_{churn} sequence in the stochastic sampling of EDM (Karras et al., 2022) framework.

F MORE DISCUSSION ON AUGMENTED RODE SAMPLING

F.1 RODE SAMPLING FOR DDIM

Recall that

$$\begin{aligned} \text{SDE: } dz &= [f(z, t) - g(t)^2 \nabla_z \log p_t(z)] dt + g(t) d\bar{w} \\ \text{RODE: } dz &= [f(z, t) - \omega(t)^2 \eta(t)^2 \nabla_z \log p_t(z)] dt + \omega(t) \eta(t) h(t) dt \end{aligned} \quad (47)$$

The insight to sample η_{t_i} from some certain probability density functions (PDFs) is inspired by sampling the scatter directions from PDFs at each hit point in ray tracing (Glassner, 1989).

Assumption 1. Consider a fixed timestep schedule $T = t_0 > t_1 > \dots > t_n = 0$. The SDE-based and RODE-based sampler share the same discretization grid and numerical integrator, thus the discretization error is the same for both schemes. We ignore the common discretization error in the subsequent analysis.

Assumption 2. The timestep interval Δt in the discrete timestep schedule is sufficiently small.

Assumption 3. For all k , η_{t_k} is independent of ϵ_{t_k} , and η_{t_k} s are independent across steps. All constants appearing below are finite, and do **not** blow up as $\Delta t \rightarrow 0$.

The exact solution of the RODE in Eq. 47 writes

$$z_{t_{i+1}} = z_{t_i} + \int_{t_i}^{t_{i+1}} [f(z_t, t) - \omega(t)^2 \eta(t)^2 \nabla_z \log p_t(z_t) + \omega(t) \eta(t) e(t)] dt, \quad (48)$$

which is approximated by

$$z_{t_{i+1}} \approx z_{t_i} + [f(z_{t_i}, t_i) - \omega_{t_i}^2 \eta_{t_i}^2 s_\theta(z_{t_i}, t_i) + \omega_{t_i} \eta_{t_i} \epsilon_{t_i}] (t_{i+1} - t_i), \quad (49)$$

where $s_\theta(z, t)$ is the score model (Song & Ermon, 2019).

The RODE sampling adopts anisotropic and non-diffusive noise, which lacks the properties of a Wiener process (e.g., uncorrelated increments and isotropic diffusion), suggesting poorer coverage of the ϵ_{t_i} space. The noise term is confined to the direction of ϵ_{t_i} and randomly scaled by η_{t_i} for each time step t_i , resembling a randomized perturbation along fixed paths rather than a true stochastic diffusion.

After switching to $\eta(t)$ as the stochastic process in sampling, the drift term $[f(z, t) - g(t)^2 \nabla_z \log p_t(z)] dt$ in the original SDE sampling also becomes random. If we assume that, reusing the SDE-trained score network $s_\theta(z, t) \approx \nabla_z \log p_t(z_t)$ within RODE sampling does **not** increase score-estimation error (this assumption will be relaxed in Prop. 6), then

- On the one hand, we argue that stochastic DDIM (Song et al., 2020) can naturally accommodate this randomness without special handling, since the learned score function (Song et al., 2021) is robust to the noise schedule σ_{t_i} (Song et al., 2021). That is, any error introduced at the current step is automatically corrected by the model in the subsequent step;
- On the other hand, this randomness mitigates the problem of reduced diversity caused by restricted process noise, which is empirically validated in Sec. 5.2 of the main paper.

We will ignore this term in the following discussion, and focus on the impact of the noise term whose stochastic is fundamentally changed. Under the above assumptions, the DDIM step can be written as

$$z_{t_{i+1}} \approx \Psi_{t_i}(z_{t_i}) + \omega_{t_i} \eta_{t_i} \epsilon_{t_i}, \quad (50)$$

where $\Psi_{t_i}(\cdot)$ is a deterministic map.

Assumption 4. Each deterministic map $\Psi_{t_i}(\cdot)$ is L_{t_i} -Lipschitz in z . Then the forward-propagation is upper bound by

$$M_{t_{i+1}} \triangleq \prod_{j=i+1}^{n-1} L_{t_j} \quad (0 \leq j \leq n-1), M_{t_n} \triangleq 1. \quad (51)$$

F.2 VARIANCE ANALYSIS

Note that $\eta(t) \in (0, 1]$ for stochastic DDIM sampling (Song et al., 2020). If $\text{Var}(\eta(t))$ exists, it can be further guaranteed that, $v_t \triangleq \text{Var}(\eta(t)) < \frac{1}{4}$ according to Popoviciu’s inequality on variances (Bencze et al., 2010), thus finite.

The contribution of the noise term in Eq. 48 is

$$\int_{t_0}^{t_n} \omega(t) \eta(t) e(t) dt. \quad (52)$$

At time step t_i in the discrete analog, the variance it brings writes

$$\text{Var}_{\text{RODE}} = \text{Var}(\omega_{t_i} \eta_{t_i} \epsilon_{t_i} \Delta t) = \omega_{t_i}^2 (\Delta t)^2 \text{Var}(\eta_{t_i}) \epsilon_{t_i} \epsilon_{t_i}^\top. \quad (53)$$

$\epsilon_{t_i} \epsilon_{t_i}^\top$ is rank-1 since it has only one linearly independent column, unlike the full-rank variance in Var_{SDE} (Eq. 42). Intuitively, this variance is anisotropic and rank-deficient, with smaller magnitude due to $\text{Var}(\eta_{t_i}) < 1$ and $(\Delta t)^2$, which results in lower variance in directions orthogonal to ϵ_{t_i} , providing lower variance and reducing sample diversity. Prop. 4 describes this property:

Proposition 4. (Prop. 1 in the main paper) *The variance of our RODE sampling is strictly less than that of standard SDE sampling in Loewner order (Horn & Johnson, 2012).*

Proof. Consider the trace of the variances of both sampling paradigm, which reflect the overall dispersion from the mean.

The trace of the standard SDE sampling’s variance (Eq. 42) at times step t_i writes

$$\text{Tr}(\text{Var}_{\text{SDE}}) = d \omega_{t_i}^2 \eta^2 \Delta t, \quad (54)$$

where $d = c \times h \times w$ denotes the dimension of ϵ_{t_i} or I .

The trace of our RODE sampling’s variance (Eq. 53) writes

$$\text{Tr}(\text{Var}_{\text{RODE}}) = \omega_{t_i}^2 \text{Var}(\eta_{t_i}) (\Delta t)^2 \|\epsilon_{t_i}\|^2, \quad (55)$$

since its eigenvalues are $\|\epsilon_{t_i}\|^2$ (with multiplicity 1) and 0 (with multiplicity $(n-1)$).

Eq. 54 scales with Δt , while Eq. 55 scales with $(\Delta t)^2$. As $\Delta t \rightarrow 0$, $\text{Tr}(\text{Var}_{\text{RODE}}) \leq \text{Tr}(\text{Var}_{\text{SDE}})$ since all the other items are finite.

More precisely, consider the Loewner order (Horn & Johnson, 2012) of these two variances. Specifically, for

$$\begin{aligned} \Delta \text{Var} &= \text{Var}_{\text{SDE}} - \text{Var}_{\text{RODE}} \\ &= \omega_{t_i}^2 \eta^2 \Delta t I - \omega_{t_i}^2 \text{Var}(\eta_{t_i}) (\Delta t)^2 \epsilon_{t_i} \epsilon_{t_i}^\top, \end{aligned} \quad (56)$$

In directions orthogonal to ϵ_{t_i} , the $\epsilon_{t_i} \epsilon_{t_i}^\top$ term contributes 0, so the eigenvalues of ΔVar in these directions are $\omega_{t_i}^2 \eta^2 \Delta t > 0$; In the direction of ϵ_{t_i} , the eigenvalue

$$\omega_{t_i}^2 \eta^2 \Delta t - \omega_{t_i}^2 \text{Var}(\eta_{t_i}) (\Delta t)^2 \|\epsilon_{t_i}\|^2 > 0 \quad (57)$$

1782 since $\eta^2 > \text{Var}(\eta_{t_i}) \|\epsilon_{t_i}\|^2 \Delta t$ as $\Delta t \rightarrow 0$ (this is because $\eta > 0$ for SDE sampling).

1783 Therefore, ΔVar is positive definite, making $\text{Var}_{\text{RODE}} < \text{Var}_{\text{SDE}}$ in Loewner order. \square

1784 The lower variance makes the distribution derived from our RODE sampling more concentrated
1785 than that from the standard SDE sampling. More precisely, under some sub-Gaussian assumptions
1786 (detailed in Appx. F.3), for constants $C, c > 0$, the vector concentration yields
1787

$$1788 \Pr \left(\|\Delta_{t_n}^{(\bullet)}\| \geq C\sqrt{V_{\bullet}} + u \right) \leq 2 \exp \left(-c \frac{u^2}{V_{\bullet}} \right), \quad (58)$$

1789 where $\bullet = \{\text{SDE}, \text{RODE}\}$, $\Delta_{t_n}^{(\bullet)} = z_{t_n}^{(\bullet)} - \mathbb{E}[z_{t_n}^{(\bullet)}]$ is the deviation of samples around their mean.
1790 Thus, both schemes display exponential tail decay and strong concentration around their means.

1791 F.3 ERROR ANALYSIS

1792 F.3.1 DISTRIBUTION DISTANCE

1793 As a novel sampling paradigm, we need to explore whether our RODE sampling can approximate
1794 the true data distribution p_{data} as the standard SDE sampling (Song et al., 2021) does. To do this, we
1795 consider the Wasserstein-1 distance (Kolouri et al., 2017) $W_1(\cdot, \cdot)$ between the distribution $p^{(S)} \triangleq$
1796 $p_{t_n}^{(S)}$ and $p^{(R)} \triangleq p_{t_n}^{(R)}$ induced by SDE and RODE sampling, respectively, where $p_t^{(\bullet)}$ is the marginal
1797 distribution at timestep t .

1800 **Assumption 5.** $p^{(S)}$ and $p^{(R)}$ have finite first moments.

1801 **Assumption 6.** The score model $s_{\theta}(z, t)$ is well-trained, so that the score-estimation error does **not**
1802 dominate the stochastic error budget. Precisely, the contribution from the model error at timestep t_i
1803 is upper-bounded by a sub-Gaussian perturbation of scale γ_{t_i} . Perturbations at different timesteps
1804 are assumed independent.

1805 We perform local linearization for SDE and RODE steps to propagate point-wise error:

$$1806 z_{t_{i+1}}^S = \Psi_{t_i}(z_{t_i}^S) + \omega_{t_i} \eta \epsilon_{t_i} + \delta_{t_i}^{\text{score}, S}, \quad (59)$$

$$1807 z_{t_{i+1}}^R = \Psi_{t_i}(z_{t_i}^R) + \omega_{t_i} \eta_{t_i} \epsilon_{t_i} + \delta_{t_i}^{\text{score}, R}, \quad (60)$$

1808 where $\delta_{t_i}^{\text{score}, \bullet}$ is the equivalent perturbation induced by the score-estimation error.

1809 **Proposition 5.** (Prop. 2 in the main paper) *The Wasserstein-1 distance of distributions induced by
1810 the two paradigms is bounded:*

$$1811 W_1(p^{(S)}, p^{(R)}) \leq \sum_{k=0}^{n-1} M_{t_{k+1}} (\omega_{t_k} \cdot C_{\epsilon, t_k} \cdot \mathbb{E}[\Delta \eta_{t_k}]) + \mathcal{O} \left(\sum_{k=0}^{n-1} M_{t_{k+1}} \gamma_{t_k} \right), \quad (61)$$

1812 where $C_{\epsilon, t_k} \triangleq \mathbb{E}[\|\epsilon_{t_k}\|]$, and $\Delta \eta_{t_k} \triangleq |\eta_{t_k} - \eta| \leq 1$.

1813 **Proof.** Consider the difference of samples.

1814 Define the point-wise trajectory error as $\Delta_{t_i} \triangleq z_{t_i}^R - z_{t_i}^S$. A first-order Taylor expansion yields

$$1815 \Delta_{t_{i+1}} = \Psi_{t_i}(z_{t_i}^R) - \Psi_{t_i}(z_{t_i}^S) + (\omega_{t_i} \eta_{t_i} - \omega_{t_i} \eta) \epsilon_{t_i} + (\delta_{t_i}^{\text{score}, R} - \delta_{t_i}^{\text{score}, S}), \quad (62)$$

$$1816 = J_{t_i} \Delta_{t_i} + r_{t_i}(\Delta_{t_i}) + \omega_{t_i} (\eta_{t_i} - \eta) \epsilon_{t_i} + \Delta_{t_i}^{\text{score}}$$

1817 where J_{t_i} is the Jacobian of $\Psi_{t_i}(\cdot)$ evaluated at the reference point in the local linearization, $r_{t_i}(\Delta_{t_i})$
1818 is the higher-order residual term, and $\Delta_{t_i}^{\text{score}} \triangleq \delta_{t_i}^{\text{score}, R} - \delta_{t_i}^{\text{score}, S}$, $\|\Delta_{t_i}^{\text{score}}\| \leq 2\gamma_{t_i}$.

1819 **Assumption 7.** The Jacobian J_{t_i} is bounded, and $\|J_{t_i}\| \leq L_{t_i}$.

1820 **Assumption 8.** The mapping $\Psi_{t_i}(\cdot)$ is twice continuously differentiable with bounded second
1821 derivatives. Thus, there exists a constant $C_{t_i} > 0$ such that $\|r_{t_i}(\Delta_{t_i})\| \leq C_{t_i} \|\Delta_{t_i}\|^2$.

1836 Neglecting the quadratic residual, the main components of the point-wise error writes
 1837

$$1838 \quad \|\Delta_{t_{i+1}}\| \leq L_{t_i} \|\Delta_{t_i}\| + \omega_{t_i} \cdot |\eta_{t_i} - \eta| \cdot \|\epsilon_{t_i}\| + \|\Delta_{t_i}^{\text{score}}\|. \quad (63)$$

1839 To make the bound more interpretable, we separate the contribution of the η -bias term $\mathbb{E}[\Delta\eta_{t_k}]$ from
 1840 the remaining score-estimation error γ_{t_k} . Starting from $\Delta_{t_0} = 0$ (i.e., shared initial noise in SDE
 1841 and RODE sampling), the difference of samples writes
 1842

$$1843 \quad \mathbb{E}\|\Delta_{t_n}\| \leq \sum_{k=0}^{n-1} M_{t_{k+1}} (\omega_{t_k} \cdot C_{\epsilon, t_k} \cdot \mathbb{E}[\Delta\eta_{t_k}] + 2\gamma_{t_k}) + (\text{higher-order residual}) \\
 1844 \quad \leq \sum_{k=0}^{n-1} M_{t_{k+1}} (\omega_{t_k} \cdot C_{\epsilon, t_k} \cdot \mathbb{E}[\Delta\eta_{t_k}]) + \mathcal{O}\left(\sum_{k=0}^{n-1} M_{t_{k+1}} \gamma_{t_k}\right), \quad (64)$$

1849 **The first term corresponds to the main contribution arising from the bias of the scalar process**
 1850 **$\eta(t)$, while the second term collects all residual errors, including the score-estimation error and**
 1851 **higher-order discretization terms.** Hence, the discrepancy between samples generated from the two
 1852 paradigms can be controlled, **and the η -bias is the dominant source of discrepancy between RODE**
 1853 **and SDE trajectories (Assump. 6).**

1854 Let $\Pi(p, q)$ be the set of coupling distributions with marginals p and q . Couple $p^{(S)}$ and $p^{(R)}$ on the
 1855 same probability space (i.e., using the same initial noise and process noise realization $\{\epsilon_{t_i}\}_{i=0}^{n-1}$).
 1856 Note that

$$1857 \quad W_1(p^{(S)}, p^{(R)}) = \inf_{\pi \in \Pi} \mathbb{E}_{\pi, X \sim p^{(S)}, Y \sim p^{(R)}} \|X - Y\| \\
 1858 \quad \leq \mathbb{E}_{\pi_0, z_{t_n}^{(S)} \sim p^{(S)}, z_{t_n}^{(R)} \sim p^{(R)}} \|z_{t_n}^{(S)} - z_{t_n}^{(R)}\| \quad (65) \\
 1859 \quad = \mathbb{E}\|\Delta_{t_n}\|,$$

1862 where π_0 is a specific coupling distribution. Thus the distance of distributions induced by the two
 1863 paradigms is also controllable. \square
 1864

1865 F.3.2 SCORE-ESTIMATION ERROR

1866 Furthermore, we relax the assumption made in Appx. F.1 by explicitly analysing the score-estimation
 1867 error that arises when RODE sampling reuses the score network $s_\theta(\cdot)$ originally trained for SDE
 1868 sampling. Intuitively,
 1869

- 1870 1. $s_\theta(\cdot)$ is trained with a time-invariant η , whereas RODE sampling employs a time-varying
 1871 η_{t_i} . This mismatch may push the latents into regions where the training data is sparse and
 1872 $s_\theta(\cdot)$ is therefore under-regularized, leading to larger score-estimation error;
 1873
- 1874 2. The standard SDE sampling adopts a Wiener process to explore all directions in the ϵ_{t_i}
 1875 space uniformly, which shows *error contraction effect* (Appx. D.4.3). However, the diffu-
 1876 sive nature of $d\bar{w}$ is absent in RODE sampling, and the directional noise lacks the ability to
 1877 provide the same error-correcting benefits as the SDE's diffusive noise, potentially allow-
 1878 ing errors to accumulate in underrepresented directions. Apart from restricted perturbation
 1879 direction, the noise term in our RODE sampling scales with dt rather than \sqrt{dt} in standard
 1880 SDE, which may limits its ability to correct accumulated error;
- 1881 3. **The $s_\theta(\cdot)$ approximates marginal distributions and is independent of sampling trajectories,**
 1882 **and RODE sampling only alters the trajectories but **not** the input distribution, therefore**
 1883 **reusing $s_\theta(\cdot)$ in RODE sampling does not introduce an additional error order.**

1884 More precisely, let the point-wise score-estimation error at timestep t be

$$1885 \quad \mathcal{E}(z, t) = \|s_\theta(z, t) - s(z, t)\|, \quad (66)$$

1887 where $s(z, t)$ is the ground-truth score function.

1888 **Assumption 9.** $\mathcal{E}(\cdot, t)$ is $L_{\mathcal{E}, t}$ -Lipschitz. Otherwise, we can resort to arguments in [Birrell \(2025\)](#),
 1889 where a local Lipschitz condition is imposed to handle the unbounded cases.

Proposition 6. (Prop. 3 in the main paper) *The difference of score-estimation error writes*

$$|\mathbb{E}_{p_t^{(R)}} \mathcal{E}(z, t) - \mathbb{E}_{p_t^{(S)}} \mathcal{E}(z, t)| \leq L_{\mathcal{E}, t} \cdot W_1(p_t^{(R)}, p_t^{(S)}), \quad (67)$$

where the constant $L_{\mathcal{E}, t} > 0$.

Proof. The Kantorovich–Rubinstein duality (Villani et al., 2008) of Wasserstein-1 distances gives

$$\begin{aligned} |\mathbb{E}_{p_t^{(R)}} \mathcal{E}(z, t) - \mathbb{E}_{p_t^{(S)}} \mathcal{E}(z, t)| &\leq L_{\mathcal{E}, t} \cdot \sup_{\text{Lip}(\mathcal{E}_0) \leq 1} \left| \mathbb{E}_{p_t^{(R)}} \mathcal{E}_0(z, t) - \mathbb{E}_{p_t^{(S)}} \mathcal{E}_0(z, t) \right| \\ &= L_{\mathcal{E}, t} \cdot W_1(p_t^{(R)}, p_t^{(S)}), \end{aligned} \quad (68)$$

where $\text{Lip}(\mathcal{E}_0) \leq 1$ denotes the set of 1-Lipschitz functions \mathcal{E}_0 . \square

Consequently, if the deviations $\mathbb{E}[\Delta\eta_{t_k}]$ s are large, RODE sampling may increase the expected score estimation error relative to SDE sampling. Conversely, if $\mathbb{E}[\Delta\eta_{t_k}]$ s are small, (e.g., when the fixed η in SDE sampling equals to the $\mathbb{E}[\Delta\eta_{t_k}]$ in RODE sampling), the above increase is small.

Summary. Our RODE sampling provides lower variance due to reduced randomness, but higher error due to biased exploration and time-varying η_{t_i} . Empirically, RODE sampling can still generate reasonable samples under larger errors (Sec. 5.2 in the main paper), which we attribute to:

1. The training process of DMs, such as the training of score functions (Song & Ermon, 2019) is intrinsically robust to noise variations (Song & Ermon, 2019; Song et al., 2021; Chen et al., 2025b);
2. DMs are robust to sampling variations, and adjusting noise schedules $\sigma(t)$ (indirectly influenced via $\eta(t)$ in our scenario) can shift sample characteristics (Karras et al., 2022);
3. Smaller noise magnitude limits dispersion, but can not prevent the “drift term” reaching the data manifold, thus ensuring plausibility.

Together, these factors justify the practical validity and stability of using RODE sampling for exploration.

F.3.3 STRONGER PROPOSITION

Assumption 10. *The score-estimation error for SDE sampling decreases as $o_N(1)$ as the sample size $N \rightarrow \infty$. Note that the analogous error term may remain at the order of $\mathcal{O}_N(1)$ for RODE sampling.*

The right-hand side in Prop. 5 may not vanish as $\Delta t \rightarrow 0$, instead, it typically converges to a finite nonzero constant. Consequently, the right-hand side of Proposition 3 does not vanish either. Therefore, Proposition 3 only establishes a stability-type guarantee that the score-estimation error of RODE sampling does not blow up, rather than demonstrating convergence to the zero-error limit achieved by SDE sampling.

Actually, under stronger assumptions, one can show mathematically that $p^{(R)}$ and $p^{(S)}$ are sufficiently close.

Assumption 11. $\mathbb{E}[\Delta\eta_{t_k}] = \mathcal{O}(\Delta t^\nu)$, $\nu > 1$.

Proposition 7. *The right-hand side in Prop. 5 vanishes, therefore:*

1. $p^{(R)}$ and $p^{(S)}$ are sufficiently close;
2. Prop. 6 demonstrate that the score-estimation error of RODE sampling converges to the zero-error limit as SDE sampling does.

Proof. The first term of the right-hand side in Prop. 5 vanishes since

$$\sum_{k=0}^{n-1} \mathbb{E}[\Delta\eta_{t_k}] = n \cdot \mathcal{O}(\Delta t^\nu) = \frac{T}{\Delta t} \mathcal{O}(\Delta t^\nu) = T \cdot \mathcal{O}(\Delta t^{\nu-1}) \rightarrow 0, \quad (69)$$

while the second term also vanishes as $N \rightarrow \infty$. \square

F.3.4 COMPARISON BETWEEN SDE SAMPLING AND RODE SAMPLING

Pseudo-Code. We provide pseudo-code in Alg. 1 and Alg. 2 to clearly demonstrate the behavior of SDE sampling and RODE sampling when generating a batch of samples.

Algorithm 1 SDE Sampling for a Batch

Input:

- B : Batch size
- n : Inference step
- $\epsilon_\theta(\cdot, \cdot)$: Noise prediction network
- $\{t_i\}_{i=0}^{n-1}$: Timestep schedule ($t_0 > t_1 \cdots > t_n$)
- $\{z_{t_0}^{(k)}\}_{k=0}^{B-1}$: A batch of initial noise

```

1: for  $i = 0$  to  $(n - 1)$  do
2:   for  $k = 0$  to  $(B - 1)$  do
3:      $\epsilon_{\text{pred}}^{(k)} \leftarrow \epsilon_\theta(z_{t_i}^{(k)}, t_i)$ 
4:      $\epsilon_{\text{process}}^{(t_i, k)} \sim \mathcal{N}(0, I)$  ◁ Sample process noise on the fly
5:      $z_{t_{i+1}}^{(k)} \leftarrow \mathcal{G}(z_{t_i}^{(k)}, \epsilon_{\text{pred}}^{(k)}, \epsilon_{\text{process}}^{(t_i, k)}, \{t_i, t_{i+1}\})$ 
6:   end for
7: end for
8: return  $\{z_{t_n}^{(k)}\}_{k=0}^{B-1}$ 

```

Process Noise. We provide an example to clarify the distinction of process noise adopted within SDE sampling and RODE sampling. Consider two samples A and B within a **single** experiment, which adopt $\{\epsilon_{t_i}^{(A)}\}_{i=0}^{n-1}$ and $\{\epsilon_{t_i}^{(B)}\}_{i=0}^{n-1}$ as process noise sequences, respectively. For SDE sampling, $\{\epsilon_{t_i}^{(A)}\}_{i=0}^{n-1}$ and $\{\epsilon_{t_i}^{(B)}\}_{i=0}^{n-1}$ may differ, while they are identical for RODE sampling.

Comparison of Distribution. For a shared realization of process noise sequence $\{\epsilon_{t_i}\}_{i=0}^{n-1}$, we argue that $p^{(S)}$ and $p^{(R)}$ overlap, but are not in an inclusion relationship:

- For time-invariant $\eta(t)$, perturbations in RODE sampling are smaller per step, exploring a less dispersed region than SDE due to $\Delta t < \sqrt{\Delta t}$ for small Δt ;
- For time-variant $\eta(t)$, RODE sampling could venture into regions beyond SDE’s reach due to varying perturbation amplitudes.

As discussed in Appx. F.3.2, $p^{(S)}$ and $p^{(R)}$ are close when the time-invariant $\eta(t)$ in SDE sampling is equal to the $\mathbb{E}[\eta(t)]$ in RODE sampling.

F.4 AUGMENTED RODE SAMPLING

We take m independent sets of process noise $\{\{\epsilon_{t_i}^{(j)}\}_{i=0}^{n-1}\}_{j=0}^{m-1}$ for our *augmented RODE sampling*. Note that this does **not** change the marginal distribution of RODE sampling.

Proposition 8. As m increases, $\{\epsilon_{t_i}^{(j)}\}_{j=0}^{m-1}$ spans \mathbb{R}^d and shows approximate isotropy due to the independence of the process noise.

Proof. The empirical covariance of the noise terms in Eq. 53

$$\frac{1}{m} \sum_{j=0}^{m-1} [\omega_{t_i} \eta_{t_i}^{(j)} \epsilon_{t_i}^{(j)} \Delta t] [\omega_{t_i} \eta_{t_i}^{(j)} \epsilon_{t_i}^{(j)} \Delta t]^\top \quad (70)$$

approximates

$$\omega_{t_i}^2 \eta_{t_i}^2 (\Delta t)^2 I_d \quad (71)$$

Algorithm 2 RODE Sampling for a Batch**Input:**

- B : Batch size
- n : Inference step
- $\epsilon_\theta(\cdot, \cdot)$: Noise prediction network
- $\{t_i\}_{i=0}^{n-1}$: Timestep schedule ($t_0 > t_1 \cdots > t_n$)
- $\{z_{t_0}^{(k)}\}_{k=0}^{B-1}$: A batch of initial noise
- $\{\epsilon_{\text{process}}^{(t_i)}\}_{i=0}^{n-1}$: A sequence of process noise, $\epsilon_{\text{process}}^{(t_i)} \sim \mathcal{N}(0, I)$ ◁ Pre-sample

```

1: for  $i = 0$  to  $(n - 1)$  do
2:   for  $k = 0$  to  $(B - 1)$  do
3:      $\epsilon_{\text{pred}}^{(k)} \leftarrow \epsilon_\theta(z_{t_i}^{(k)}, t_i)$ 
4:      $z_{t_{i+1}}^{(k)} \leftarrow \mathcal{G}(z_{t_i}^{(k)}, \epsilon_{\text{pred}}^{(k)}, \epsilon_{\text{process}}^{(t_i)}, \{t_i, t_{i+1}\})$ 
5:   end for
6: end for
7: return  $\{z_{t_n}^{(k)}\}_{k=0}^{B-1}$ 

```

since $\sum_{j=0}^{m-1} \epsilon_{t_i}^{(j)} (\epsilon_{t_i}^{(j)})^\top \approx mI_d$ holds for large m when $\epsilon_{t_i}^{(j)}$ are *i.i.d.* random variables from $\mathcal{N}(0, I_d)$. Thus, it provides approximately isotropic exploration to all directions similar to a Wiener process, *i.e.*, the directionality of the process noise becomes less significant as m increases. \square

Note that the noise term is still scales by Δt in RODE sampling (Eq. 49) rather than $\sqrt{\Delta t}$ in SDE sampling (Eq. 40).

F.5 COMPARISON WITH PREVIOUS WORKS

Previously studies have examined the low-dimensional representations of diffusion sampling trajectories. For example, PAS (Wang et al., 2024) employs a 4D coordinate to specify the direction of a single denoising step for PF-ODE sampling (Song et al., 2021) (deterministic sampling). Differently, our approach adopts RODE sampling (based on stochastic sampling), where a single parameter η_{t_i} suffices to determine a step. Based on these observations, we conjecture that, stochastic samplers can likewise admit low-dimensional parameterizations. We leave this for future work.

Previously, Flow-GRPO (Liu et al., 2025b) introduces stochasticity to flow matching (Albergo & Vanden-Eijnden, 2022; Lipman et al., 2022; Liu et al., 2022) by converting ODE sampling to SDE sampling. Our proposed RODE sampling offers an alternative solution to introduce randomness to ODE-based approaches, which achieves lower variance than SDE-based methods, and can be used as a stable and controllable exploration policy for RL training.

G MORE DISCUSSION ON OUR MCTS

G.1 OVERVIEW

Fig. 2 in the main paper presents an overview of our MCTS with online updated Beta policies. Specifically,

1. *Selection Phase*: Traverse the tree nodes with depth in range $[0, n']$, and select a not fully expanded node with maximum value-based UCB;
2. *Expansion Phase*: Sample an a_{t_i} from the Beta policy of the selected node s_{t_i} and create a new node $s_{t_{i+1}}$, then compute the intermediate reward $R(s_{t_i}, a_{t_i}, s_{t_{i+1}})$ with any latent reward policy. Note that the Beta policy of the expanded node is not initialized, which can also be seen as initialized as a uniform distribution. Additionally, for max-reward modeling, the pseudo-final samples obtained during the intermediate reward calculation are used to update the best trajectory if achieving higher rewards than the best-known trajectory;

- 2052 3. *Simulation Phase*: Sample an $\eta_{t_{i+1}} \sim \mathcal{U}([0, 1])$ for the expanded node as the initial action,
 2053 and initialize its unimodal Beta policy via mode re-parameterization. Then, perform a
 2054 stochastic simulation with uniformly sampled actions sequentially until reaching a terminal
 2055 state, while computing intermediate rewards along the roll-out. If a trajectory superior to
 2056 the best-known one is found, the Beta policies from the current node to the root node in the
 2057 MCTS tree are updated via mode re-parameterization;
- 2058 4. *Backpropagation Phase*: Propagate the obtained merged reward of the terminal state back
 2059 through the visited nodes with max-value policy;
- 2060 5. Loop the above steps until the NFEs have elapsed.

2062 G.2 SUPPLEMENTARY TO MODE RE-PARAMETERIZATION

2063
 2064 **Flaws of Directly Updating the Shape Parameters.** Directly update the shape parameters is
 2065 straightforward, but hard to balance the exploration and exploitation of the Beta policy, since α and
 2066 β are prior counts of successes and failures rather than probabilities of sampling certain actions.

2067 **The Concentration Controller ζ .** Consider a unimodal Beta distribution $\text{Beta}(\alpha, \beta)$, where
 2068 $\alpha, \beta > 1$, which are re-parameterized as $\alpha = 1 + e^a, \beta = 1 + e^b$, respectively. The mode and
 2069 the variance of it are $\rho = \frac{\alpha - 1}{\alpha + \beta - 2}$ and $\sigma^2 = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$, respectively.

2070 We introduce a hyper-parameter $\zeta > 0$ for concentration control, and set $a = \ln \zeta, b =$
 2071 $\ln \zeta + \ln \frac{1 - \rho}{\rho}$. Actually, for a fixed $\rho \in (0, 1)$, the changes in σ^2 as ζ increases are non-monotonic.

2072 However, we empirically found that, setting ζ appropriately (we adopt $\zeta = 3$) can control the con-
 2073 centration of the unimodal Beta distribution to some extent (Appx. N.2). Nevertheless, the derivative
 2074 of σ^2 with respect to $t = \ln \zeta$ changes sign at points denpent on the specific value of $\rho \in (0, 1)$.
 2075 This means that, we can **not** directly determine a reasonable range for ζ , so its selection relies on
 2076 hyper-parameter grid search. We leave designing a better re-parameterization for the unimodal Beta
 2077 distribution with its mode for future research.

2081 G.3 SUPPLEMENTARY TO INITIALIZATION AND ONLINE UPDATE

2082
 2083 When a trajectory superior to the best-known one is found, the Beta policies from the current node
 2084 to the root node within the MCTS tree are softly updated. Specifically, assume an action η^* sampled
 2085 from the Beta policy with mode of ρ yields better performance, the mode is updated to

$$2086 \rho' \leftarrow \begin{cases} \eta^*, & \text{if } |\eta^* - \rho| \leq \kappa \\ \rho + \kappa \cdot \text{sgn}(\eta^* - \rho), & \text{otherwise} \end{cases}, \quad \rho' \leftarrow \text{Clip}_\varepsilon^{1-\varepsilon}(\rho'), \quad (72)$$

2087 where ε is a small positive value to ensure unimodality, and $\kappa > 0$ is the update step size, which
 2088 limits the amplitude amplitude to ensure policy stability (Lee et al., 2020).

2091 G.4 EFFICIENCY OPTIMIZATION

2092
 2093 **Batch Processing.** The MCTS processes for all samples are conducted in a batch-wise manner,
 2094 sharing a single search tree. Nodes at depth i ($0 \leq i \leq n$) can store at most a batch of states
 2095 (latents). When expanding the k -th sample in a batch, if the node selected for expansion has a
 2096 child node with an empty k -th state, we prioritize populating this child node rather than creating
 2097 a new one. This reduces the number of tree nodes. Meanwhile, we sort each node’s children in
 2098 non-ascending order of their weight (the number of filled states), and place emptier nodes earlier in
 2099 the children list for faster traversal.

2100 **VRAM Management.** All node states expanded during MCTS process are stored on the GPU. To
 2101 efficiently manage GPU memory (VRAM), we utilize a global LRU cache (Fricker et al., 2012) to
 2102 maintain the states of tree nodes. We set a global limit χ for the number of GPU-resident states, *i.e.*,
 2103 the capacity of the cache. During expansion, if the number of states on the GPU is not less than χ ,
 2104 we offload the least recently used state in the cache to the CPU, and move the expanded one to the
 2105 GPU. When accessing a node’s state, if it resides on the CPU, we repeat the aforementioned cache-
 maintenance steps to ensure the returned latents on GPU. We set $\chi = 1000$ in our experiments.

2106 H EXPERIMENTAL SETTINGS, EVALUATION DIMENSIONS AND REWARD 2107 FUNCTIONS 2108

2109 H.1 IMPLEMENTATION 2110

2111 We modified the the HuggingFace diffusers (von Platen et al., 2022) library for implementation.
2112

2113 H.2 SETTINGS 2114

2115 **Environment.** All experiments are conducted on a *Ubuntu 20.04.6 LTS* system with $128 \times$ *Intel(R)*
2116 *Xeon(R) Gold 6430* CPU and $1 \times$ *RTX 6000 Ada Generation GPU* with *48G VRAM*.
2117

2118 **Text-to-Image Models.** We adopt the following pre-trained models as our base models for align-
2119 ment:

- 2120 • *stabilityai/sd-turbo* (Sauer et al., 2024b) ¹;
- 2121 • *compvis/stable-diffusion-v-1-4-original* (Sauer et al., 2024b) ²;
- 2122 • *stabilityai/stable-diffusion-xl-base-1.0* (Podell et al., 2023) ³;
- 2123 • *PixArt-alpha/PixArt-XL-2-1024-MS* (Chen et al., 2023b) ⁴.
2124
2125

2126 We use the full-precision versions of the models, because we empirically found that, the fp16 version
2127 can cause VAE (Kingma et al., 2013; 2021) precision issues, making the decoded images completely
2128 black (Madebyollin, 2024).

2129 We do **not** adopt the current SoTA SD 3 (Esser et al., 2024), SD 3.5 (Esser et al., 2024),
2130 FLUX.1 (Labs et al., 2025; Labs, 2024) or Hunyuan DiT (Li et al., 2024b) models, as they are
2131 trained under flow-matching objectives (Albergo & Vanden-Eijnden, 2022; Lipman et al., 2022; Liu
2132 et al., 2022), and are incompatible with the DDIM scheduler (Song et al., 2020). We leave general-
2133 ization to flow-matching schedulers for future work.
2134

2135 **Hyper-Parameters.** We adopt the following hyper-parameters unless specified.

- 2136 • The global random seed is set to 42;
- 2137 • The seeds for initial noise is $\{0, 1, 2, \dots\}$;
- 2138 • For SD-Turbo and SD v1.4, the sample resolution is 512×512 , with a CFG (Ho & Salimans,
2139 2022) scale of 6.5; For SDXL base and PixArt- α , the sample resolution is 1024×1024 ,
2140 with a CFG scale of 4.5. The *negative_prompt* is “low quality, blurry, ugly,
2141 oversaturated” for all models. Particularly, for SDXL base and PixArt- α , the
2142 *prompt_2* and *negative_prompt_2* is set to an empty string;
- 2143 • For MCTS, the hyper-paramter that balance exploration and exploitation $\lambda_{\text{exploitation}} = 2.0$,
2144 the selection depth limit is $n' = \left\lceil \frac{4n}{5} \right\rceil$, and simulation actions are uniformly sampled from
2145 $[0, 1]$;
- 2146 • For the Beta policies, the update step size $\kappa = 0.1$, and the clamp epsilon to $\varepsilon = 10^{-8}$.
2147 The concentration control scalar is set to $\zeta = 3$ unless specially specified.
2148
2149
2150

2151 H.3 METRIC 2152

2153 H.3.1 MPD 2154

2155 We use LPIPS (Zhang et al., 2018) to evaluate the perceptual similarity between two images,
2156 where values closer to 0 indicate greater visual similarity. We calculate the mean pairwise distance

2157 ¹<https://huggingface.co/stabilityai/sd-turbo>

2158 ²<https://huggingface.co/CompVis/stable-diffusion-v-1-4-original>

2159 ³<https://huggingface.co/stabilityai/stable-diffusion-xl-base-1.0>

⁴<https://huggingface.co/PixArt-alpha/PixArt-XL-2-1024-MS>

(MPD) (Kim et al., 2025) to assess the diversity of a set of samples $\{x_i\}_{i=0}^{N-1}$:

$$\text{MPD} = \frac{1}{N(N-1)} \sum_{i=0}^{N-2} \sum_{j=i+1}^{N-1} \text{LPIPS}(x_i, x_j), \quad (73)$$

where larger MPD means higher diversity. The MPD of generation results per prompt is calculated and averaged across all prompts to gauge the diversity of a specific inference-time scaling paradigm.

H.3.2 PARAMETER EFFICIENCY

The number of parameters adopted for the scaling process is defined as

$$\#\text{Param.} = m \cdot n \cdot \dim(\text{action}), \quad (74)$$

where $\dim(\text{action})$ is the dimension of the adopted action. Specifically, in our experiments,

- If scaling with η_{t_i} , then $d_\eta \triangleq \dim(\eta_{t_i}) = 1$;
- If scaling with ϵ_{t_i} , then $d_\epsilon \triangleq \dim(\epsilon_{t_i})$, which is equal to $4 \times 64 \times 64$ for SD v1.4, or $4 \times 128 \times 128$ for SDXL base and PixArt- α .

The *Parameter Efficiency* (PE, \uparrow) quantifies how available parameters are efficiently used for the scaling process, which writes

$$\text{PE} = \frac{\text{reward} - \text{baseline}}{\#\text{Param.}}, \quad (75)$$

where *reward* is the maximum achieved reward, and *baseline* is the baseline reward for normalization.

H.3.3 EVALUATION DIMENSIONS

We consider the following dimensions to evaluate our MCTS:

- *Best Reward*: The averaged highest reward achieved, which is the maximum final reward for cumulative-reward modeling, while it is the maximum encountered reward for max-reward modeling. It is rounded to 4 decimal places;
- *NFE-dynamics*: The averaged NFE consumed for computing the transition dynamics. It is rounded to 2 decimal places, so do *NFE-intermediate* and *NFE-final*;
- *NFE-intermediate*: The averaged NFE consumed for calculating intermediate rewards;
- *NFE-final*: The averaged NFE consumed for calculating final rewards. Note that the NFEs needed for calculating intermediate rewards are included;
- *Time Cost*: The averaged wall-clock time required per sample to obtain *Best Reward*. It is rounded to the nearest integer.

For each prompt, we run our MCTS with m groups of process noise and take the optimal results (the samples with highest best rewards), and average the evaluation metrics like *Best Reward* across these results. If multiple samples with the same best rewards exist, the one with the smallest NFEs and the shortest time will be retained. Particularly, since our implement MCTS in batch (Appx. G.4), the time cost per sample is influenced by other samples in the batch. We record the wall-time consumed to achieve the best reward for each sample individually, and report the averaged time cost.

H.3.4 REWARD FUNCTIONS

We adopt the following models or metrics as reward functions to align DMs with:

- **HPS v2** (Human Preference Score v2) (Wu et al., 2023b;a)⁵: Quantifies the human aesthetics preference of images. Higher scores indicate better aesthetics. Particularly, we report 100 times of the original reward;

⁵<https://github.com/tgxs002/HPSv2>

- **IR** (ImageReward) (Xu et al., 2023a)⁶: Evaluates the human preference of images. Higher scores indicate higher preference;
- **PS** (PickScore) (Kirstain et al., 2023)⁷: Predicts human preference of images. Higher scores indicate higher preference. Note that we evaluate the PS in sample-wise manner, instead of aggregating the samples into a batch and computing in a single pass;
- **CLIP Score** (Radford et al., 2021; Hessel et al., 2021)⁸: Evaluates the semantic alignment between the text descriptions and their generated images. Higher scores indicate higher alignment. We adopt *laion/CLIP-ViT-H-14-laion2B-s32B-b79K*⁹ for feature extraction;
- **Compressibility Reward** (CR): The file size obtained by compressing the image with the standard JPEG compression (Raid et al., 2014) from Python’s Pillow library (Appx. H.3.6), with a *quality* parameter of 30 and enabling optimization;
- **LAPV** (Laplacian Variance, LAPV) (Memon et al., 2015; GeeksforGeeks, 2024): The variance of the Laplacian matrix (GeeksforGeeks, 2024) of an image, which is an **unbounded** value that quantifies the sharpness of images. The larger, the sharper;
- **CCR** (Color Channel Reward, CCR) (Eyring et al., 2024): Measures how much an RGB image’s overall hue leans towards red (R), green (G) or blue (B). For an RGB image $X_{3 \times h \times w}$ with height h and width w , let the pixel value at coordinate (i, j) in channel c ($0 \leq c \leq 2$) be $X_{c,i,j} \in [0, 1]$ (normalized). If the target channel is c , and the other two channels are \bar{c}_1 and \bar{c}_2 , then the degree to which X leans towards channel c can be quantified as

$$\text{CCR}(c) = \sum_{i,j} X_{c,i,j} - \sum_{i,j} X_{\bar{c}_1,i,j} \sum_{i,j} X_{\bar{c}_2,i,j}. \quad (76)$$

H.3.5 REWARD PRE-PROCESSING

Some rewards should be pre-processed to meet the non-negativity requirement for MCTS values. We first clip the rewards in a specific *range* (if specified), then perform normalization by dividing them by a *scaling factor*, and finally add a *bias scalar*. These constants are listed in Tab. 11.

H.3.6 METRIC COMPUTATION

We utilize the following third-party libraries for metric computation:

- *Human Preference Score v2* (1.2.0) (Wu et al., 2023a)¹⁰;
- *pytorch-fid* (0.3.0) (Yu et al., 2021)¹¹;
- *lpips* (0.1.4) (Zhang et al., 2018)¹²;
- *opencv-python* (4.11.0.86) (OpenCV, 2024)¹³;
- *open-clip-torch* (2.23.0) (Ilharco et al., 2021)¹⁴;
- *Pillow* (2.23.0)¹⁵.

Table 11: **Reward pre-processing.**

Reward Type	Clip Range	Scaling Factor	Bias Scalar
HPS v2	/	1.0	0.0
PS	/	1.0	0.0
IR	$[-2, 2]$	1.0	2.0
CLIP Score	/	1.0	0.0
CR	/	1.0	3.0
LAPV	/	1.0	0.0
CCR	/	1.0	2.0

⁶<https://github.com/THUDM/ImageReward>

⁷<https://github.com/yuvalkirstain/PickScore>

⁸https://github.com/mlfoundations/open_clip

⁹<https://huggingface.co/laion/CLIP-ViT-H-14-laion2B-s32B-b79K>

¹⁰<https://github.com/tgxs002/HPSv2>

¹¹<https://github.com/mseitzer/pytorch-fid>

¹²<https://github.com/richzhang/PerceptualSimilarity>

¹³<https://github.com/opencv/opencv-python>

¹⁴https://github.com/mlfoundations/open_clip

¹⁵<https://github.com/python-pillow/Pillow>

H.4 OVERVIEW OF EXPERIMENTAL COMPOSITION

Sec. 5.2 in the main paper and Appx. I provide the empirical foundation for our RODE-based sampling; Sec. 5.4 and Appx. J align DMs with aesthetics, and makes comparison between different inference-time scaling paradigms; Sec. 5.5 and Appx. K align DMs with semantics, and benchmark our method against Z-Sampling (Bai et al., 2024); Sec. 5.6 and Appx. L show our method’s ability to align with composite rewards; Appx. M show the generalization of our method across various latent reward policies; Sec. 5.7, Appx. N and Appx. O presents ablation studies and applications. Appx. P presents failure cases.

I EXPERIMENTS: SUPPLEMENTARY TO RODE SAMPLING

I.1 SUPPLEMENTARY TO TRAJECTORY SIMULATION

In Fig. 1 in Sec. 5.2 of the main paper, we simulate the trajectories produced by three sampling paradigms using an Ornstein-Uhlenbeck (OU) process (Karatzas & Shreve, 2012)

$$dx_t = -\mu x_t dt + \sigma_t dw_t \tag{77}$$

with the discretized step

$$x_{t_{i+1}} = x_{t_i} - \mu x_{t_i} \Delta t + \sigma_{t_i} \sqrt{\Delta t} \epsilon_{t_i} \tag{78}$$

for n times sequentially.

Specifically,

- For PF-ODE sampling (Song et al., 2021), we fix $\sigma_{t_i} = 0$ for $i = 0, 1, \dots, n - 1$;
- For SDE sampling (Song et al., 2021), we fix $\sigma_{t_i} = 1$, and re-sample the process noise ϵ_{t_i} at every timestep t_i ;
- For RODE sampling, we pre-sample $\{\epsilon_{t_i}\}_{i=0}^{n-1}$ and re-use them across trajectories. We sample σ_{t_i} from $\mathcal{U}([0, 1])$ at each timestep t_i .

We set $n = 1000$, $\mu = 1.5$, $\Delta t = \frac{1}{n}$. The background heatmap depicting the evolution of data distributions is generated by applying a Gaussian filter with $\sigma = 3$ to the 2D histograms of 5k SDE trajectories.

I.2 SUPPLEMENTARY TO SAMPLE DIVERSITY

Settings. We randomly sample 100 prompts from the HPD v2 dataset (Wu et al., 2023a). For each prompt, $N = 5$ images are generated under the conditions in Tab. 13 using 2-step SD-Turbo (Sauer et al., 2024a) and SD v1.4 (Rombach et al., 2022) with 15, 20, and 25 steps, respectively. we adopt the $m = 3$ sets of process noise for Aug. RODE sampling. The MPD of generation results per prompt is calculated and averaged across all prompts to gauge the diversity.

Table 12: Seeds used for each paradigm.

Paradigm	Initial Noise	Process Noise
a.0	{0, 1, 2, 3, 4, 5}	/
a.1	{0, 1, 2, 3, 4, 5}	$\mathcal{U}([3072, 4095])$
a.2	{0, 1, 2, 3, 4, 5}	{3072, 3073, ...}
b	{0, 0, 0, 0, 0}	$\mathcal{U}([3072, 4095])$
c.0 / c.1	{0, 0, 0, 0, 0}	{3072, 3073, ...} {4096, 4097, ...} {5120, 5121, ...}

Policies of Different Inference-Time Scaling Paradigms. The initial noise and process variance are easy to fix by keeping their values constant. However, it is infeasible to fix the process noise by setting all $\{\epsilon_{t_i}\}_{i=0}^{n-1}$ to a constant vector ϵ , because this time-invariant perturbation will cause large accumulative error in SDE sampling, leading to significant divergence of true trajectories and producing completely noisy images in multi-step inference scenarios. For fair comparison, we set seeds for each paradigm as in Tab. 12. Particularly, for *paradigm c*, we run with $m = 3$ sets of process noise, and compute MPD for each set of the process noise **respectively**, which is averaged across the $m = 3$ sets as the result.

Table 13: **Policies and MPDs (\uparrow) of different inference-time scaling paradigms.** Cells with a red / orange / yellow background: the best / second-best / third-best performance for each model (column), respectively.

Paradigm	Index	Sampler	Components			SD-Turbo		SD v1.4	
			$\{z_{t_0}^{(j)}\}_{j=0}^{N-1}$	$\{\epsilon_{t_i}\}_{i=0}^{n-1}$	$\{\eta_{t_i}\}_{i=0}^{n-1}$	2-step	15-step	20-step	25-step
a) initial-noise	a.0)	ODE	random	\times	$\eta_{t_i} = 0$	0.5867	0.7542	0.7584	0.7594
	a.1)	SDE	random	random	$\eta_{t_i} = 1.0$	0.6437	0.7653	0.7617	0.7650
	a.2)	Aug. RODE	random	m sets	$\eta_{t_i} = 1.0$	0.6200	0.6861	0.6887	0.6914
b) process-noise		SDE	fixed	random	$\eta_{t_i} = 1.0$	0.3823	0.6848	0.6939	0.6946
c) process-noise -variance	c.0)	Aug. RODE	fixed	m sets	$\eta_{t_i} \sim \mathcal{U}([0, 1])$	0.1628	0.4717	0.4852	0.4899
	c.1)	Aug. RODE	fixed	m sets	$\eta_{t_i} \sim \mathcal{U}([0.5, 1])$	0.1295	0.3556	0.3649	0.3686

Results and Analysis. We provide analysis of results in Tab. 13 as follows:

- Under random initial noise, all sampling paradigms exhibit high sample diversity. Specifically, the initial noise (*paradigm a*) predominately determines the sampling diversity, especially in few-step inference (e.g., $n = 2$ and $n = 15$);
- The process noise (*paradigm b*) impacts sampling diversity more significantly than the process variance (*paradigm c*);
- Though we fixed the process noise and randomize the variance, *paradigm c* shows no mode collapse (Barsha & Eberle, 2025) but mid-to-high diversity, especially when n is large;
- As the number of inference steps increases, diversity rises across all paradigms. Particularly, in few-step inference, *paradigm b* and *paradigm c* show much lower sampling diversity than *paradigm a*. However, as n increases, the diversity in *paradigm b* and *paradigm c* significantly rises, with *paradigm b* even surpassing *paradigm a.2*. Our Aug. RODE sampling (*paradigm c.0*) even achieves an MPD of nearly 0.5 when generated with 25-step SD v1.4 (Rombach et al., 2022). It indicates that, the sampling diversity in multi-step inference can be boosted by fully utilizing the process noise and process variance at each denoising step to perturb the denoising trajectories;
- Interestingly, for Aug. RODE sampling, adopting a broader range $[0, 1]$ for η_{t_i} (*paradigm c.0*) yields higher diversity than a larger scale $[0.5, 1]$ (*paradigm c.1*). On the one hand, this inspires us to adopt $\eta_{t_i} \sim \mathcal{U}([0, 1])$ in our experiments. On the other hand, this highlights that, the exploration can be easily adjusted by altering the perturbation magnitude (i.e., the range of η_{t_i});
- Our Aug. RODE sampling (*paradigm c*) partially integrates the ability of SDE sampling (*paradigm b*);
- The empirical observations above suggest that, our method is more suitable for multi-step inference scenarios due to broader exploration.

In summary, the augmented RODE sampling successfully balances the controllability of ODE sampling and the exploration of SDE sampling, serving as a qualified and efficient sampler for denoising trajectory search.

I.3 SUPPLEMENTARY TO EMPIRICAL DISTRIBUTION

Motivation. To deploy RODE sampling in practical generation tasks, we should additionally verify that, its empirical distribution does not excessively deviate from the real data distribution even with larger errors (Prop. 5, Prop. 6).

Settings. The MS-COCO 5k validation set (Lin et al., 2014) comprises 5,000 images, each paired with 5 or 6 captions for a total of 25,010 prompts. Selecting SD v1.4 (Rombach et al., 2022) as the base model, we generate $N = 1$ image per prompt with various inference steps using ODE (Song et al., 2020), SDE (Song et al., 2021) and Aug. RODE sampling, respectively. We adopt $m = 3$ sets of process noise for SDE sampling and Aug. RODE sampling, the zero-shot FID is computed for each set of the process noise **respectively**, and is averaged across the $m = 3$ sets as the result.

2376 **Results and Analysis.** Results in Tab. 3 demonstrate that,
2377

- 2378 1. Aug. RODE sampling attains FIDs that lie between those of ODE and SDE sampling,
2379 indicating that the novel paradigm using $\eta_{t_i} s$ as random variables does not significantly
2380 shift the distribution. We attribute this to the close resemblance of the drift term of Aug.
2381 RODE to that of the original SDE (Appx. F.3). Additionally, we empirically find it hard to
2382 distinguish RODE-sampled results from standard SDE/ODE-sampled ones (e.g., Fig. 1);
- 2383 2. The ordering $\text{FID}_{\text{SDE}} < \text{FID}_{\text{Aug. RODE}} < \text{FID}_{\text{ODE}}$ highlights the forgetting effect induced
2384 by injecting process noise within stochastic sampling (Appx. D.4.3, Appx. F.3);
- 2385 3. This empirically corroborate the observations in Song et al. (2021) and Karras et al. (2022)
2386 that, given sufficient NFE, introducing process noise improves the sample quality.
2387

2388 I.4 VISUALIZATION OF SAMPLING TRAJECTORIES

2389 **Visualization.** We detail the establishment of the polar plot for visualizing DM’s high-
2390 dimensional denoising trajectories. Consider a series of concentric semicircles $\{\Gamma_i\}_{i=1}^n$ in polar
2391 coordinates, with an angular range of $[0, \pi]$, and radii of $\frac{1}{n}, \frac{2}{n}, \dots, 1$, respectively. The origin, as
2392 the 0-th level semi-circle, represents the initial noise z_{t_0} which is the denoising start point. The i -th
2393 ($1 \leq i \leq n$) semi-circle denotes the data manifold of the noise level at time step t_i . By linear map-
2394 ping each $\eta_{t_i} \in [0, 1]$ to $\eta'_{t_i} \in [0, \pi]$, the denoising trajectory parameterized by the action sequence
2395 $\eta_{t_0, \dots, t_{n-1}}$ can be visualized as follows: starting from the origin, at each i -th step ($0 \leq i \leq n - 1$),
2396 moving from i -th level manifold along the η'_{t_i} -angle direction to $(i+1)$ -th level manifold, eventually
2397 reaching n -th level and obtaining a final latent z_{t_n} .
2398
2399

2400 **Settings.** We visualize the sampling trajectories and the corresponding qualitative results
2401 for the prompt “A symmetrical oil painting of two waterfalls in a dense
2402 forest.” in Fig. 4. It demonstrate that,
2403

- 2404 1. The trajectory of deterministic DDIM (Song et al., 2020) consistently follows the 0-angle
2405 direction, the one of DDPM (Ho et al., 2020) consistently follows the π -angle direction,
2406 while our RODE sampling navigates time-variant directions;
- 2407 2. Our RODE sampling explores various paths from the origin to the n -th level data manifold
2408 (the outmost semi-circle), seeking for higher rewards.

2409 A comprehensive study on the relationship between action sequences and sampling results is left for
2410 future work.
2411

2412 **Comparison with Previous Works.** Previous works (Zhou et al., 2024a; Chen et al., 2024; Wang
2413 et al., 2024; Wang & Vastola, 2023) that low-dimensional visualizations of diffusion sampling tra-
2414 jectories are confined to deterministic PF-ODE (Song et al., 2021) sampling. Our approach offers a
2415 novel 2D perspective that renders stochastic sampling trajectories more readily and intuitively inter-
2416 pretable than prior representations, advancing the visualization of stochastic sampling paradigm.
2417

2418 J EXPERIMENT: SUPPLEMENTARY TO ALIGNING WITH AESTHETICS

2419 J.1 SETTINGS

2420 **The Dataset, Model and Reward Function.** We randomly sample 50 prompts from the HPD
2421 v2 dataset (Wu et al., 2023a), and sample $N = 2$ images for each prompt. We align 15-step SD
2422 v1.4 (Rombach et al., 2022) to HPS v2 (Wu et al., 2023a), aiming at enhancing the sample aesthetics.
2423
2424
2425

2426 **NFE Budgets.** The limit of NFE dynamics is set to 999.
2427

2428 **Methods.** We contend that, existing inference-time scaling approaches with ϵ_{t_i} actions (Oshima
2429 et al., 2025; He et al., 2025; Liu et al., 2025a; Mao et al., 2025; Singhal et al., 2025) all boil down to
uniformly sampling $\epsilon_{t_i} \in \mathbb{R}^{c \times h \times w}$ and applying diverse and elaborate optimizations. We strip away

additional embellishments, and essentially evaluate ϵ_{t_i} -action scaling in isolation. Specifically, we instantiate beam search (BS) (Fernandes et al., 2025; Oshima et al., 2025), global search (GS, the BS with a single beam) and MCTS Coulom (2006); Browne et al. (2012) as base planners, and equip each with either ϵ_{t_i} or η_{t_i} actions. Particularly, for ϵ_{t_i} -action methods, we adopt $m = 3$ random sets of process noise generated using random seeds sampled from $\mathcal{U}([3072, 4095])$, $\mathcal{U}([4096, 5119])$, and $\mathcal{U}([5120, 6143])$, respectively, and retain the highest rewards (similar to $m = 3$ fixed sets of process noise that η_{t_i} -action methods do). In summary, we make comparison between:

- Vanilla DDIM (Song et al., 2020) sampling as baseline;
- Best-of- N (Ma et al., 2025), which is **not** directly comparable as it is z_{t_0} -action and falls outside our setting, but is included here as a widely used baseline. To respect the NFE budget, we require $n \cdot N \leq 999$, yielding $N \leq 66$; thus we set $N = 66$. Note that this N refers to the number of candidate initial noise in best-of- N , which is **not** related to the N appearing in “sample $N = 2$ images for each prompt”;
- Inference-step scaling paradigm (Nichol & Dhariwal, 2021) using DDPM (Ho et al., 2020) and DDIM (Song et al., 2020). Note that the DDPM sampling with $m = 3$ sets of process noise can also be seen as an ϵ_{t_i} -action scaling. For both Best-of- N and inference-step scaling, the reported *Time Cost* refers exclusively to the sampling time and does **not** include the time required for reward computation. This is because *Time Cost* is defined as the average wall-clock time per sample needed to obtain the *Best Reward*, and these two methods must generate a fixed number of samples regardless of when the best rewards are found;
- Beam search (Fernandes et al., 2025; Oshima et al., 2025) with $\epsilon_{t_i}/\eta_{t_i}$ actions that employs the conventional max-final-reward modeling and *immediate-ddim* as latent reward policy. It should be noted that, our definitions of B and K might slightly differ from those in the original paper (Oshima et al., 2025). Specifically, for a fixed initial noise, we maintain B beams at every depth of the search tree. Each beam stochastically spawns K candidate nodes, of which only the top- B performers are retained for expansion at the next layer. Particularly, at depth 1, we directly initialize with $B \cdot K$ candidates in one shot. Besides, to respect the NFE budget, we enforce $n \cdot B \cdot K \leq 999$, yielding $B \cdot K \leq 66$. Following Oshima et al. (2025), we restrict B and K to powers of two, which further tightens this constraint to $B \cdot K$. we exhaustively enumerate all (B, K) pairs that satisfies $B \cdot K = 64$, and report the best-performing configuration. We also report the performance of greedy search (GS), which corresponds to $(B, K) = (1, 64)$;
- Our MCTS with $\epsilon_{t_i}/\eta_{t_i}$ actions. Particularly, for MCTS-eps (ϵ_{t_i} action), we adopt settings identical to MCTS-eta (η_{t_i} action), except the expansion policy for ϵ_{t_i} being changed as described above.

Notes. It is worth mentioning that, since deterministic DDIM (Song et al., 2020) does not employ process noise, a direct comparison with DDPM (Ho et al., 2020) and other methods — which use m sets of process noise — is unfair if strictly speaking. Nevertheless, we retain this evaluation to:

- Highlight the performance gains afforded by leveraging multiple sets of process noise;
- Treat the performance of the deterministic DDIM as the baseline from which relative improvements are computed.

J.2 RESULTS AND ANALYSIS

Comparison on the Scaling Process. Fig. 5 shows that, increasing inference steps yields limited sample quality enhancement, and can even degrade performance with too many steps, which is consistent with the observations in Nichol & Dhariwal (2021); Li et al. (2023). In contrast, our MCTS achieves higher rewards and faster convergence, which underscores the superiority of allocating the NFE budgets across multiple denoising trajectories rather than concentrating it on a single one. BS methods are omitted because they only yield samples at the maximum depth.

Generalization to Other Aesthetic Scores. We also evaluate the optimization process of HPS v2 (Wu et al., 2023a) with PickScore (PS) (Kirstain et al., 2023) and ImageReward (IR) (Xu et al.,

2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537

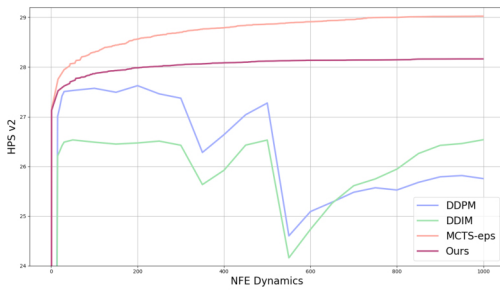


Figure 5: Comparison on the optimization process of HPS v2 between different methods.

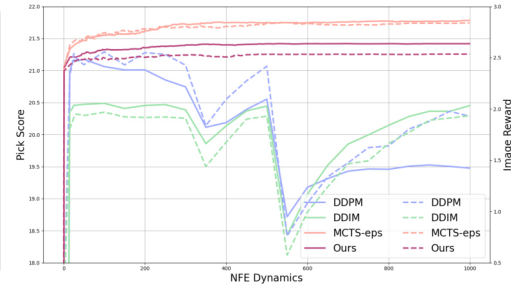


Figure 6: Comparison of PS and IR between our RODE-based scaling and inference-step paradigm. Solid/dashed lines for PS/IR, respectively.

2023a). Results in Fig. 6 show improvements in two additional human aesthetics preference predictor that weren’t directly optimized, indicating that our MCTS truly enhances image aesthetics rather than just over-optimizing HPS v2. We omit the evaluation of Aesthetic Score (Romain Beaumont, 2022) here and in Appx. O.2, since it gives discrete scores for each image.

Details of quantitative results for inference-step scaling and our method are provided in Tab. 15 and Tab. 16, respectively.

Quantitative Results and Parameter Efficiency. Quantitative results are shown in Tab. 14, where relative improvement is computed based on 15-step DDIM. It demonstrates that,

1. GS, BS and our MCTS all outperform the scaled DDIM and DDPM;
2. Ours (MCTS-eta) surpasses traditional inference-step scaling. Note that scaling DDPM is essentially an ϵ_{t_i} -action method;
3. When scaling with η_{t_i} , ours outperform both GS and BS, highlighting MCTS’s ability to allocate NFEs judiciously, and Beta policies’ capacity for accurate low-dimensional search;
4. Although ϵ_{t_i} -action methods achieve the highest absolute score, its relative improvement is only $\sim 2\times$ better than η_{t_i} -action ones, but with $\sim 16k\times$ more parameters, revealing low parameter efficiency and severe parameter redundancy;
5. When naively transferring our MCTS to ϵ_{t_i} -actions, performance drops below GS and BS, underscoring the limitation of vanilla MCTS in high-dimensional continuous action spaces (Bianchi et al., 2023);
6. Ours runs significantly faster than GS and BS. This is because GS and BS expand nodes in strictly increasing depth order. In contrast, MCTS can revisit shallow but high-value nodes in later search, and these nodes spend more NFEs for expansion and simulation;
7. MCTS obviates the need to enumerate (B, K) pairs, which is more adaptive for deployment.

Qualitative Results. Samples displayed in Fig. 1 (a) are generated with text prompt “A tea kettle sits on the burner of stove.”, where samples for DDPM and DDIM are selected from the inference steps that yields the highest HPS v2. It can be observed that, trajectory search-based methods render more reasonable kettles and avoid artifacts (e.g., the fire on a cup in DDPM). Note that it is hard to distinguish samples derived by SDE sampling and RODE sampling.

More qualitative results in Fig. 7 highlight our superior visual aesthetics over baselines. The text prompts are “Yoda performing at Woodstock.”, “A dreamlike scene with a vaporwave aesthetic.”, “A manga-style illustration of a submachine gun in 2050 by Moebius and Stephan Martiniere.”, “Psytrance artwork by Lisa Frank.”, and “The image is a digital art poster sized in the

Table 16: Aesthetic scores of our method.

Metric	HPS v2	PS	IR
Reward	28.1639	21.4383	2.5350
Improvement (%)	7.45	5.27	40.23

Table 14: **Comparison between different methods when optimizing HPS v2.** “Relative Improv.”: relative improvement.

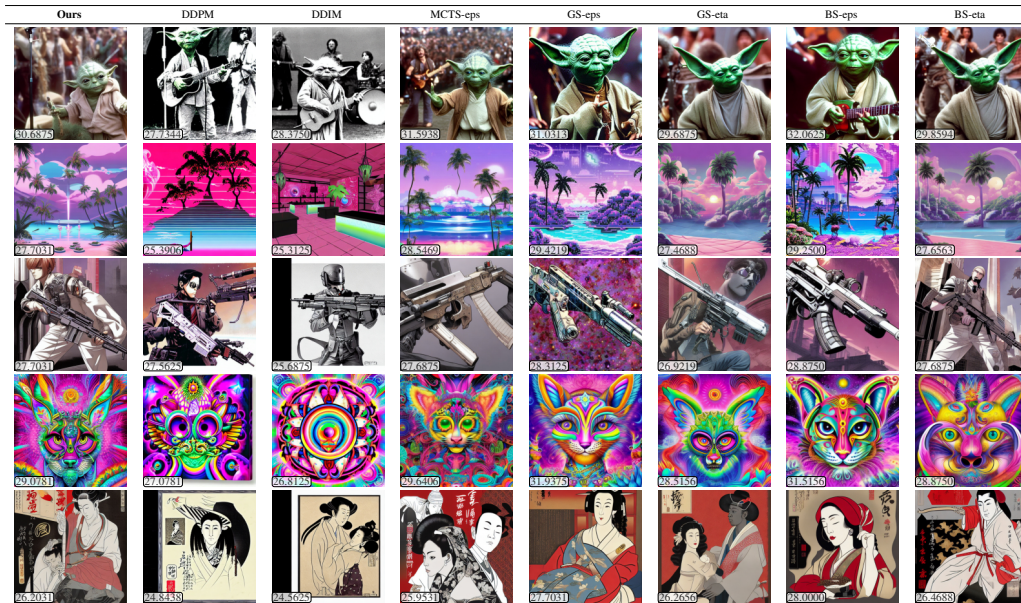
Method	Action Type	Best Reward	Relative Improv. (%)	#Param.	Param. Efficiency	Time Cost (s)
DDIM (15 steps, baseline)	/	26.2122	/	/	/	147
Best-of-N (DDIM, $N = 66$)	z_{t_0}	28.4056	8.37	$66 \cdot d_\epsilon$	2.03e-6	4480
Best-of-N (DDPM, $N = 66$)	z_{t_0}	28.4297	8.46	$66 \cdot d_\epsilon$	2.05e-6	4513
DDIM (999 steps)	ϵ_{t_i}	26.5362	1.24	/	/	6694
DDPM (200 steps)	ϵ_{t_i}	27.6256	5.39	$3 \cdot 200 \cdot d_\epsilon$	1.44e-7	1382
GS ($B = 1, K = 64$)	ϵ_{t_i}	29.7386	13.45	$3 \cdot 15 \cdot d_\epsilon$	4.78e-6	5684
BS ($B = 8, K = 8$)	ϵ_{t_i}	30.2095	15.25	$3 \cdot 15 \cdot d_\epsilon$	5.42e-6	5469
MCTS-eps	ϵ_{t_i}	29.0241	10.73	$3 \cdot 15 \cdot d_\epsilon$	3.81e-6	1437
GS ($B = 1, K = 64$)	η_{t_i}	27.8116	6.10	$3 \cdot 15 \cdot d_\eta$	0.0355	5697
BS ($B = 8, K = 8$)	η_{t_i}	28.0833	7.14	$3 \cdot 15 \cdot d_\eta$	0.0416	5273
Ours (MCTS)	η_{t_i}	28.1644	7.45	$3 \cdot 15 \cdot d_\eta$	0.0434	2736

Table 15: **Quantitative results of the inference-step scaling for SD v1.4.** HPS v2 (\uparrow), PickScore (\uparrow), and Image Reward (\uparrow), from top to bottom.

Steps	15	20	25	30	50	100	150	200
	250	300	350	400	450	500	550	600
	650	700	750	800	850	900	950	999
DDPM	27.0064	27.1961	27.4134	27.5067	27.5283	27.5727	27.4923	27.6256
	27.4605	27.3734	26.2814	26.6403	27.0405	27.2764	24.6028	25.0923
	25.2805	25.4803	25.5706	25.5252	25.6798	25.7909	25.8164	25.7555
DDIM	26.2122	26.3084	26.4150	26.4887	26.5333	26.4881	26.4502	26.4725
	26.5097	26.4264	25.6347	25.9236	26.4288	26.5312	24.1603	24.7339
	25.2594	25.6117	25.7486	25.9464	26.2603	26.4241	26.4608	26.5362
DDPM	20.9863	21.0963	21.1661	21.1489	21.1745	21.0616	21.0089	21.0098
	20.8538	20.7450	20.1144	20.1872	20.3922	20.5539	18.7192	19.1784
	19.3145	19.4295	19.4641	19.4592	19.5066	19.5255	19.5041	19.4777
DDIM	20.3647	20.3925	20.4592	20.4584	20.4728	20.4875	20.4081	20.4550
	20.4686	20.3856	19.8591	20.1367	20.3700	20.4452	18.4348	19.0692
	19.5242	19.8527	19.9955	20.1458	20.2830	20.3644	20.3636	20.4516
DDPM	2.3326	2.4108	2.5399	2.4956	2.4315	2.5571	2.4299	2.5497
	2.5331	2.4270	1.8368	2.0955	2.2742	2.4163	0.7615	1.0822
	1.3388	1.4753	1.6227	1.6395	1.8046	1.8789	1.9772	1.9300
DDIM	1.8078	1.8577	1.9151	1.9508	1.9366	1.9685	1.9234	1.9172
	1.9220	1.9097	1.4397	1.6762	1.9015	1.9294	0.5751	0.9958
	1.2412	1.4644	1.4950	1.6642	1.7719	1.8908	1.9116	1.9339

style of Utamaro Kitagawa featuring Lil Wayne.”, respectively. It can be observed that, trajectory search-based methods yield samples that are visually more appealing than those produced by inference-step scaling, whereas ours additionally avoids color over-saturation.

Figure 7: Qualitative results of optimizing HPS v2.



Note. We provide the following clarification regarding the global search capabilities of our method and absolute best rewards:

1. Given the limited expressive capacity of 1-D η -actions, η -based search behaves more like a *local search* method, and thus it is natural that its absolute best rewards may fall short of ϵ -based *global search* methods. However, experiments in Sec. 5.4 and Sec. 5.5 show that even with an extremely low-dimensional action space, our method outperforms several eps-action baselines (e.g., inference-step scaling and Z-Sampling), highlighting its high parameter efficiency;
2. As noted in Appx. Q.3, future work may explore *global-local hybrid strategies*, e.g., using ϵ -actions to identify high-reward regions, followed by η -actions for refined local search;
3. Nevertheless, the difference between the *relative improvement* achieved by η - and ϵ -based methods — when viewed through the lens of their vastly different parameter dimensionalities — reveals a broader conceptual insight: *parameter efficiency* plays a crucial role in search-based alignment.

K EXPERIMENT: SUPPLEMENTARY TO ALIGNING WITH SEMANTICS

K.1 SETTINGS

The Dataset, Model and Reward Function. We randomly sample 30 prompts from the *colors*, *descriptions* and *positional* categories in DrawBench (Saharia et al., 2022) dataset, and align 30-step SDXL base (Podell et al., 2023) to CLIP score (Radford et al., 2021; Hessel et al., 2021), aiming at enhancing semantic alignment.

NFE Budgets. The limit of NFE dynamics is set to 999.

Methods. We benchmark our method against:

- 30-step DDPM (Ho et al., 2020) and DDIM (Song et al., 2020). We do **not** adopt the scaled ones since the number of inference steps is not directly related to prompt alignment;
- Z-Sampling (Bai et al., 2024), which simultaneously enhances sample aesthetics and semantic alignment. We adopt deterministic DDIM sampling (since stochastic samplers

underperform their deterministic counterparts owing to larger approximation errors) with hyper-parameters recommended in the official implementation¹⁶ of Z-Sampling. Specifically, we apply the zig-zag operation along the entire trajectory, executing it k time(s) per denoising step for scaling ($k = 1$ in the official release of Z-Sampling), and retain the samples with the highest CLIP Score across all k s. The *back-tracking stepsize* is set to 1 for each zig-zag operation. The guidance scale for denoising and inversion are set to 5.5 and 0.0, respectively. Note that, both denoising and inversion steps are counted toward the *NFE dynamics*. To respect the NFE budget of 999, the constraint $k \cdot 3(n - 1) + 1 \leq 999$ yields $k \leq 11$.

K.2 RESULTS AND ANALYSIS

Qualitative Results. Note that, the displayed samples of Z-Sampling in Fig. 1 (d) in the main paper and below are produced with $k = 2$ (highest CLIP Score). We provide more qualitative results in Fig. 8. The text prompts are “A separate seat for one person, typically with a back and four legs.”, “A mechanical or electrical device for measuring time. ”, “An umbrella on top of a spoon.”, “A tennis racket underneath a traffic light.”, “A donut underneath a toilet.”, and “A train on top of a surfboard.”, respectively.

It can be observed that,

1. For line 1 to line 4, our method produces high-quality images that align well with the prompts, whereas baselines render the semantics to some extent with lower quality;
2. For line 5 to line 6, our samples accurately depict the positional relationship between the objects, which the baselines fail to capture.

Discussion on Quantitative Results. Tab. 6 in Sec. 5.5 of the main paper presents quantitative results, where the relative improvement is computed based on 30-step DDIM. It demonstrates that,

1. Our method achieves the most significant performance enhancement;
2. Z-Sampling yields larger gains in aesthetic scores. However, aesthetic scores and CLIP score are only weakly — and sometimes even negatively — correlated, *i.e.*, samples judged more aesthetically pleasing may receive lower CLIP scores, which is called *inconsistent preferences* in Zhang et al. (2025). Consequently, such undirected and inconsistent optimization fails to maximize the CLIP score.

Scaling Z-Sampling. As shown in Tab. 17, scaling Z-Sampling does not guarantee improved performance, and can induce quality degradation, *e.g.*, color saturation and noise. In contrast, our method exhibits robust scaling behavior. Qualitative results are shown in Fig. 9), with text prompts “A bicycle on top of a boat.” and “A red colored car.”, respectively. We attribute this to:

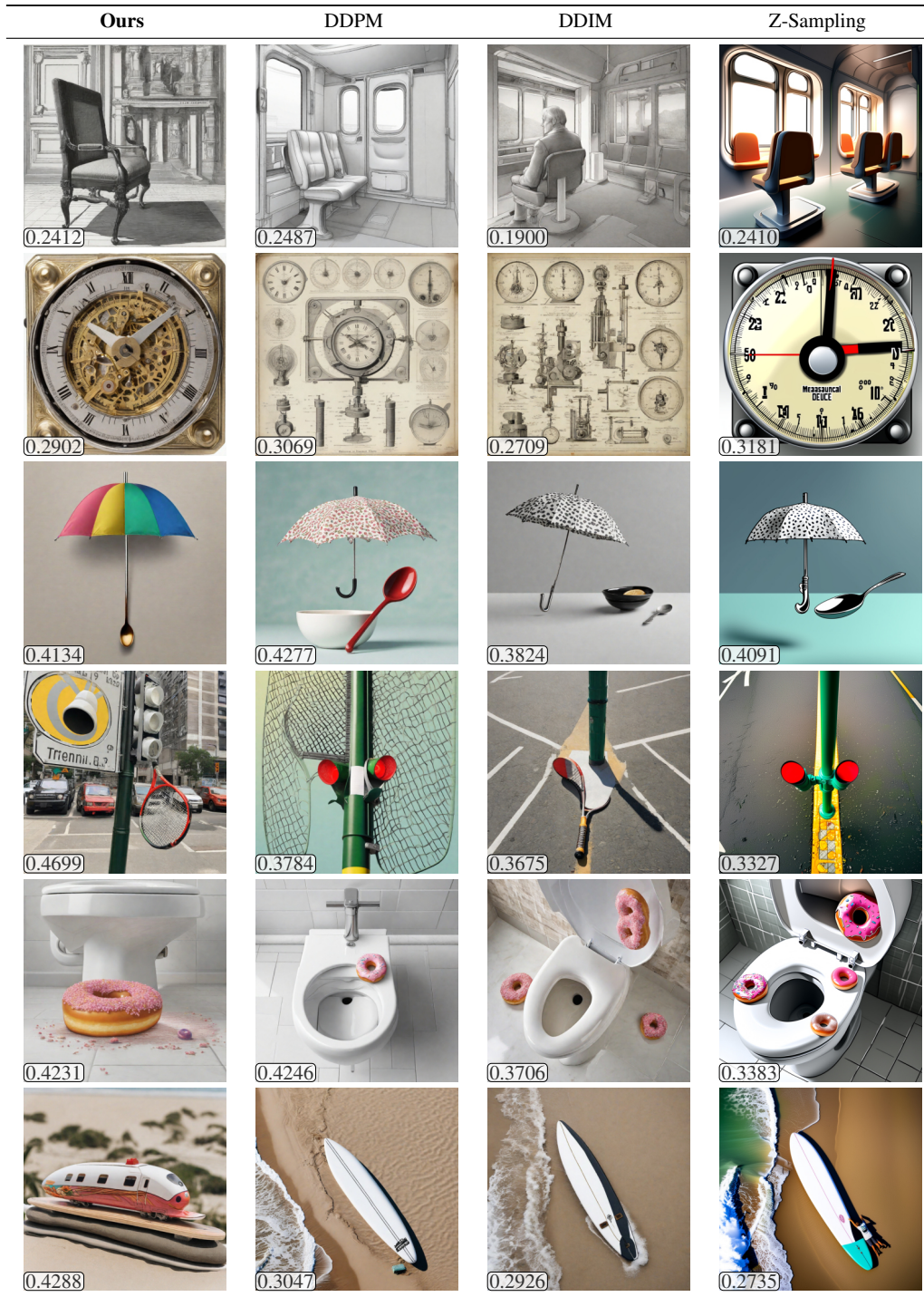
1. Information corruption introduced by repeated noising-denoising cycles (Lugmayr et al., 2022; Wang et al., 2022; Bansal et al., 2023; Karras et al., 2022);
2. Z-Sampling progressively refines the attention masks of the foreground subjects (Bai et al., 2024). An excessive number of zig-zag operations tend to blur the background or yield an overly monotonous backdrop.

Table 17: **Quantitative results of scaling Z-Sampling.** “*max*”: retain the samples with the highest CLIP Score across all k s.

k	1 7	2 8	3 9	4 10	5 11	6 max
CLIP	0.3443	0.3508	0.3500	0.3484	0.3464	0.3459
Score (↑)	0.3449	0.3418	0.3423	0.3441	0.3467	0.3676
HPS	29.0346	29.2260	29.0893	28.8760	28.7349	28.6070
v2 (↑)	28.4049	28.1979	28.0820	28.0305	27.8464	29.5635

¹⁶<https://github.com/xie-lab-ml/Zigzag-Diffusion-Sampling/>

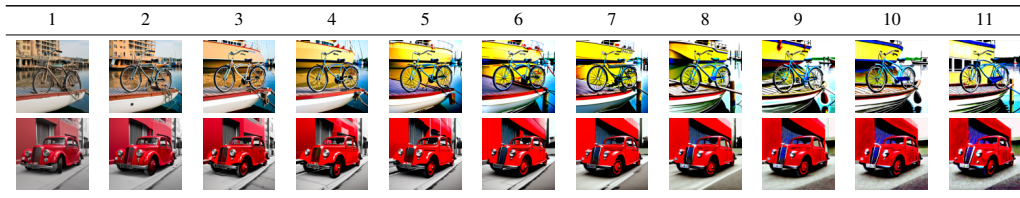
Figure 8: Qualitative results of optimizing CLIP score.



L EXPERIMENT: SUPPLEMENTARY TO ALIGNING WITH COMPOSITE REWARDS

L.1 SETTINGS

The Dataset and Model. We randomly sample 20 prompts from the HPD v2 dataset (Wu et al., 2023a), and align 50-step Pixart- α (Chen et al., 2023b) to different reward functions.

Figure 9: Qualitative results of scaling Z-sampling (various ks).

Composite Reward. The *composite reward* that combines CR and HPS v2 (Wu et al., 2023a) is defined as

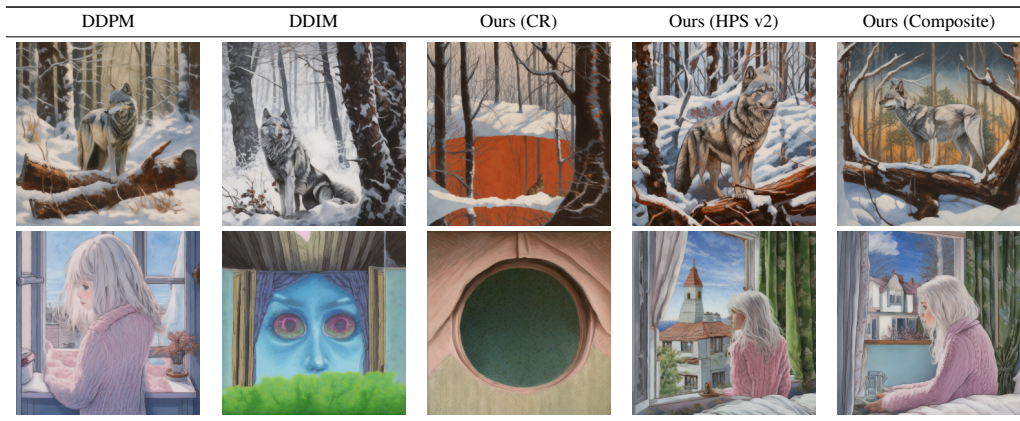
$$R_{\text{composite}} = R_{\text{compressibility}} + \lambda \cdot R_{\text{HPS.v2}}. \quad (79)$$

The following discusses the selection of hyper-parameter λ . After normalizing both compressibility reward and HPS v2 to $[0, 1]$, the former remains tightly clustered above 0.95, whereas the latter spans a much wider range. We tune λ to balance these two objectives, prioritizing samples with higher HPS v2 while still optimizing compressibility. The value of $\lambda = 0.02$ is chosen to approximate the ratio of the standard deviations of these two rewards across several random samples.

NFE Budgets. The limit of NFE dynamics is set to 999.

Methods. We conduct comparison between *Ours* (CR), *Ours* (HPS v2), and *Ours* (Composite), where *Ours* (R) denotes aligning with reward R using our method.

L.2 RESULTS AND ANALYSIS

Figure 10: Qualitative results of aligning PixArt- α with different objectives.

Qualitative results in Fig. 1 (b) in the main paper are generated with the text prompt “A beautiful blue and pink sky overlooking the beach.”. We provide more qualitative results in Fig. 10. The text prompts are “A black wolf standing on a fallen tree in a winter forest.” and “A white-haired girl in a pink sweater looks out a window in her bedroom.”, respectively. It demonstrates that,

1. When optimizing compressibility reward alone, a form of reward hacking emerges: the model collapses to introducing unreasonable smooth abnormalities to inflate the reward (e.g., Fig. 1 (b)), producing only smooth backgrounds and omits the foreground subject (e.g., Fig. 10 line 1), or even exhibit complete semantic loss (e.g., Fig. 10 line 2);
2. Optimizing the composite reward yields samples that are perceptually smoother yet retain clear semantics, validating our method’s ability to align DMs with various rewards simultaneously.

M EXPERIMENT: SUPPLEMENTARY TO GENERALIZATION TO OTHER LATENT REWARD POLICIES

M.1 SETTINGS

The Dataset, Model and Reward Function.

We randomly sample 20 prompts from the HPD v2 dataset (Wu et al., 2023a). We align 15-step SD v1.4 (Rombach et al., 2022) with HPS v2 (Wu et al., 2023a).

Latent Reward Policies.

We adopt different latent reward policies in Appx. D.2 respectively for evaluation. Particularly, the actions used for computing intermediate rewards in *sequential* policy are uniformly sampled from $[0, 1]$.

NFE Budgets. The limit of NFE dynamics is set to 500.

M.2 RESULTS AND ANALYSIS

Quantitative results in Tab. 18 demonstrate that,

1. Different latent reward policies vary in the accuracy of intermediate rewards, the quality of pseudo-final latents, and computational cost;
2. The *LA-2* and *LA-3* policies yield the 1st/2nd best performance, respectively. They take two or three steps to compute an intermediate reward, which increases the probability of obtaining better pseudo-final samples, leading to performance superior to that of the two *immediate* policies;
3. The *sequential* policy, which spends more steps calculating intermediate rewards, achieves the worst performance, as it wastes a significant amount of NFEs to calculate intermediate rewards, which should be used for exploration in other latent reward policies. This underscores the importance of efficient intermediate reward computation for effective search.

N EXPERIMENT: SUPPLEMENTARY TO ABLATIONS

N.1 SUPPLEMENTARY TO REWARD HACKING AND EFFECTS OF τ

Settings. We randomly sample 20 prompts from the HPD v2 dataset (Wu et al., 2023a), and align 50-step SD v1.4 (Rombach et al., 2022) to the Laplacian variance (LAPV) (Memon et al., 2015; GeeksforGeeks, 2024) and 3 types of color channel reward (CCR) (Eyring et al., 2024) (*CCR-R/G/B* for the R/G/B channel), respectively. The limit of NFE dynamics is set to 500.

Supplementary to Reward Hacking. Sec. 5.7 in the main paper illustrates a reward hacking phenomenon (Eisenstein et al., 2023) that, the optimization process identifies vulnerabilities in LAPV and CCRs, where noisy and blurred samples can yield higher rewards (Black et al., 2023). Exploiting this, the process leverages early-stage pseudo-final samples from our max-reward modeling to maximize rewards with unreasonable results. Luckily, this form of reward hacking can be mitigated by excluding pseudo-final samples from valid sample consideration. But this might compromise the core features of the max-reward modeling, thus potentially degrading the performance.

Results and Analysis. The displayed samples in Fig. 1 (e) in the main paper are generated with text prompts “A tall chicken standing next to a farmer.”, “A close-up hyperrealistic oil painting of a nurse fashion model with red lipstick, ginger hair, freckles, in a style mixing classicism and 80s sci-fi, set in complete darkness.”, “A symmetrical oil painting of two waterfalls in a dense forest.” and “Pippi is tethered to the international space station in her space suit amidst stars and galaxies.”, respectively. We present the following discussion:

Table 18: Performance of different latent reward policies when optimizing HPS v2.

Reward Policy	Best Reward	NFE Dynamics	NFE Inter.	NFE Final	Time Cost
<i>immediate-ddim</i>	29.2582	316.52	279.38	595.90	686
<i>immediate-score</i>	29.0293	325.88	289.27	615.15	1458
<i>LA-2</i>	29.3957	376.90	334.95	711.85	1693
<i>LA-3</i>	29.3938	326.00	290.82	616.83	1528
<i>sequential</i>	28.7070	292.07	35.90	74.55	314

1. Even though LAPV is an unbounded reward whose values span a wide range, our method achieves robust inference-time DM alignment;
2. Optimizing CCR is really challenging for our approach. Specifically, we fix the parameters of the pre-trained SD v1.4 and treat it as the environment model. Since RODE sampling can only induce limited perturbations along certain directions (determined by the process noise), our method can **not** achieve an overall color-tone effect like that in Eyring et al. (2024), as this would deviate from the original data distribution embedded in SD v1.4. However, qualitative results in Fig. 1 (e) show that, our method can sometimes obtain higher rewards in tricky ways. For example, the search process resorts to adding small components with colors closely relative to the target channel (*i.e.*, the primary color R/G/B and its secondary colors (Hunt, 2005)) for higher rewards, while maintaining reasonable coherence. This highlights the adaptability of our method in extreme cases.

N.2 ABLATIONS ON m AND ζ

Effects of m . We study the role of the hyper-parameter m in Aug. RODE sampling by limiting the number of groups of process noise. Specifically, when the full budget comprises m sets of process noise, an ablation that uses only $m' \in [1, m]$ sets is conducted by enumerating all $C_m^{m'}$ possible combinations, and reporting the best result across them. We randomly sample 20 prompts from the HPD v2 dataset (Wu et al., 2023a), and align 15-step SD v1.4 (Rombach et al., 2022) to HPS v2 (Wu et al., 2023a). Results in Tab. 19 shows that, too few groups confine exploration to limited perturbation directions, causing performance degradation. It underscores the need to augment RODE sampling.

Table 19: Ablations on m for augmented RODE sampling.

m	1	2	3 (Ours)
Best Reward	28.7344	29.1270	29.2582

Table 20: Ablations on ζ for Beta policies. Performance worse than uniform policies (29.1152 in Tab. 21) is marked with \downarrow .

ζ	1	2	3	4	5	6	7	8
	9	10	11	12	13	14	15	
Best Reward	29.1102	29.1254	29.2582	29.1410	29.1148 \downarrow	29.0547 \downarrow	29.0973 \downarrow	29.0961 \downarrow
	29.0988 \downarrow	29.0199 \downarrow	29.0371 \downarrow	29.0625 \downarrow	29.1273	29.0859 \downarrow	29.1227	

Effects of ζ . We also examined the concentration control scalar ζ in Beta policies. Specifically, we randomly sample 20 prompts from the HPD v2 dataset (Wu et al., 2023a), and align 15-step SD v1.4 (Rombach et al., 2022) to HPS v2 (Wu et al., 2023a). Tab. 20 presents the quantitative results, whose plot is shown in Fig. 11. It indicates that, $\zeta = 3$ achieves the best performance, which we attribute to an appropriate balance between exploration and exploitation. Moreover, the performance of Beta policy with an inappropriate selection of ζ may fall behind that of the naive uniform policy, especially when ζ is large. We attribute this to the Beta policy becoming trapped in local optima due to over-exploitation, whereas the uniform policy enables more thorough exploration with sufficient NFEs.

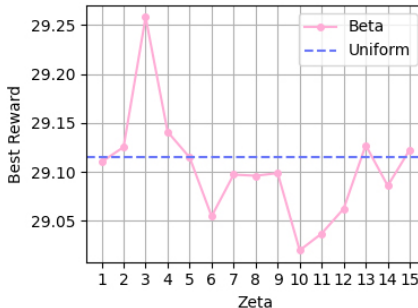


Figure 11: Ablations on ζ .

N.3 SUPPLEMENTARY TO MAIN ABLATIONS

Settings. We conduct ablations on key components of our method. Specifically, we randomly sample 20 prompts from the HPD v2 dataset (Wu et al., 2023a), and align 15-step SD v1.4 (Rombach et al., 2022) to HPS v2 (Wu et al., 2023a).

Results and Analysis. Quantitative results in Tab. 21 demonstrate that:

Table 21: **Ablations.** “Ours w/o pseudo-final”: disabling pseudo-final samples as valid samples. “Ours w/o Beta policy”: uniform policy.

Method	Best Reward	NFE Dynamics	NFE Inter.	NFE Final	Time Cost
Ours	29.2582	316.52	279.38	595.90	686
Ours w/o pseudo-final	29.0199	293.73	259.82	553.55	612
Ours w/o expansion depth limit	29.1082	298.45	259.18	557.62	631
Ours w/o online update	29.1547	325.88	287.18	613.05	1028
Ours w/o Beta policy	29.1152	297.02	264.75	561.77	528

1. Excluding pseudo-final samples from valid samples leads to the worst performance, since the samples encountered during the search are not fully leveraged. This underscores the advantage of our max-encountered-reward formulation over the conventional max-final-reward formulation;
2. In the absence of an expansion-depth cap, a disproportionate amount of NFEs is assigned to the nodes whose depth is close to n during the later search process. This is because, deeper nodes are more likely to yield higher rewards when using *immediate-ddim* policy. However, modifying η_{t_i} at late stages leads to limited alterations to latents, leaving the shallower and potentially more promising nodes under-explored;
3. Disabling the online update of Beta policy freezes the peak of the Beta distribution at its initial value, preventing the search focus from incorporating best-so-far knowledge;
4. The uniform policy underperforms Beta policy, highlighting the appropriate balance between exploitation and exploration;
5. Any component is indispensable.

O APPLICATIONS

O.1 SYNERGY WITH COMMUNITY MODULES

Our method lies in trajectory search techniques, which is plug-and-play compatible with community modules and other inference-time scaling approaches, further boosting performance. For example, Promptist (Hao et al., 2023) can be introduced to optimize the input prompts into model-preferred ones, and the initial noises for each prompt can be respectively optimized with Golden Noise (GN) (Zhou et al., 2024b) before the MCTS process.

Settings. We randomly sample 20 prompts from the HPD v2 (Wu et al., 2023a) dataset, and align the 30-step SDXL base (Podell et al., 2023) to HPS v2 (Wu et al., 2023a). We sequentially integrate our method with Promptist (Hao et al., 2023) and GN (Zhou et al., 2024b). The limit of NFE dynamics is set to 500 to highlight how optimizing the prompt and the initial noise contributes to the early search of our method.

Table 22: **Quantitative results of synergy with community modules.**

Method	HPS v2 (\uparrow)	CLIP Score (\uparrow)
Ours	28.2766	0.3564
+ Promptist	29.1418	0.4177
+ GN	29.1945	0.4181

Results and Analysis. Quantitative and qualitative results are presented in Tab. 22 and Fig. 12. The samples in Fig. 12 in the main paper is generated with text prompts:

1. Line 1: “A man smiles as he stirs his food in the pot.”, and its optimized version “a man smiling as he stirs his food in the pot by greg rutkowski, digital art, realistic painting, fantasy, very detailed, trending on artstation”;
2. Line 2: “A painting of a Persian cat dressed as a Renaissance king, standing on a skyscraper overlooking a city.”, and its optimized version “painting of a Persian cat dressed as a

Figure 12: Qualitative results of synergy with community modules.



Renaissance king, standing on a skyscraper overlooking a city, by Greg Rutkowski and Raymond Swanland, Trending on Artstation, ultra realistic digital art”;

- Line 3: “ Two vespas parked next to a light post.” and its optimized version “two vespas parked next to a light post, hyperdetailed, artstation, cgsociety, 8 k”.

Note that the reward used during search process is computed with the original prompts, whereas the metrics reported in Tab. 22 are evaluated with the optimized prompts. Results demonstrate that,

- Robust improvements persist even after prompt or initial-noise refinement;
- Our approach is orthogonal to both Promptist and GN, and these three techniques can synergistically elevate the performance of inference-time DM alignment, achieving samples with higher aesthetics.

O.2 ROBUSTNESS OF IMAGE REWARD FUNCTIONS

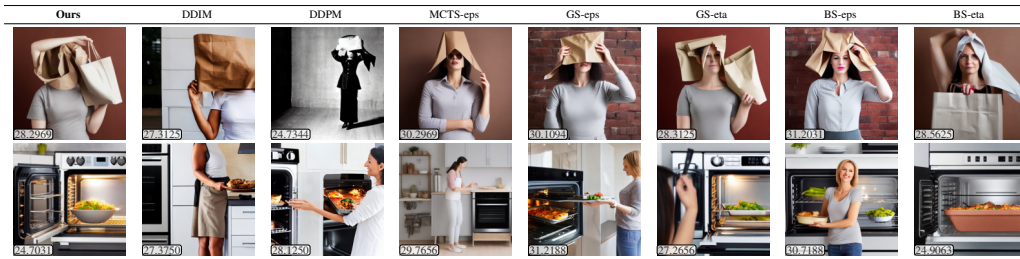
Robustness evaluation is often overlooked in existing studies on reward models for text-to-image generation (Liu et al., 2024), particularly in terms of quantitative assessment. Our RODE-based scaling method introduces limited and controllable perturbations to the denoising trajectory, resulting in a sequence of progressively similar intermediate images. This enables us to assess the robustness of the optimized (Lipschitz continuous) reward models.

Quantifying Robustness. We approximate the absolute value of the derivative of the reward change with respect to the perceptual change (measured by LPIPS (Zhang et al., 2018)) to quantify robustness. Formally, let z_{t_n} and z'_{t_n} be two consecutive intermediate samples during the optimization of a reward model ϕ . The *local robustness* of ϕ on this sample pair is measured as

$$\text{Robustness}(\phi; z_{t_n}, z'_{t_n}) = \left| \frac{\phi(z_{t_n}) - \phi(z'_{t_n})}{\text{LPIPS}(\mathcal{D}(z_{t_n}), \mathcal{D}(z'_{t_n})) + \varepsilon} \right|, \tag{80}$$

where $\mathcal{D}(\cdot)$ is the decoder of the VAE (Kingma et al., 2013; 2021) for latent DMs (Sauer et al., 2024a), and $\varepsilon > 0$ is a small constant for numerical stability. A higher robustness value indicates a less robust ϕ . The *robustness* is measured as the average of its local robustness for all initial noises and prompts.

Figure 13: Failure cases of optimizing HPS v2.

Table 23: Robustness (\downarrow) of three popular aesthetics models.

Reward Model	HPSv2	PS	IR
Robustness	1.5926	1.0674	0.6491

Settings. We randomly sample 20 prompts from the HPD v2 dataset (Wu et al., 2023a), and set $\varepsilon = 10^{-6}$. We adopt HPS v2 (Wu et al., 2023a), PS (Kirstain et al., 2023) and IR (Xu et al., 2023a) as reward models, respectively. The limit of NFE dynamics is set to 500.

Results. Tab. 23 reports the robustness of three popular aesthetic reward models, where IR exhibits the highest robustness.

O.3 QUANTITATIVE EVALUATION OF INITIAL NOISE POTENTIAL

Our method can be seen as a technique that sufficiently utilizes the potential of initial noise *i.e.*, it reveals the best achievable performance of an initial noise under given conditions and computational constraints. This can be used to quantitatively evaluate the *quality* or *generation potential* of initial noise, *e.g.*, serves as a verifier. We leave this for future work.

P FAILURE CASES

P.1 ALIGNING WITH AESTHETICS

We provide failure cases when optimizing HPS v2 (Wu et al., 2023a) (Sec. 5.4, Appx. J) in Fig. 13. The text prompts are “A portrait of a woman with a paper bag over her head.” and “A woman that is standing near an open oven.”, respectively. It can be observed that,

1. For line 1, both ours and GS-eps introduce anomalies in limbs and fingers;
2. For line 2, all three η_{t_i} -action methods underperform the baselines in score, and fails to render the semantics of “a woman” correctly.

P.2 ALIGNING WITH SEMANTICS

We provide failure cases when optimizing CLIP score (Radford et al., 2021; Hessel et al., 2021) (Sec. 5.5, Appx. K) in Fig. 14. The text prompts are “A hair drier underneath a sheep.” and “A zebra underneath a broccoli.”, respectively. It can be observed that, although our method attains higher CLIP scores, it fails to render accurate positional relationship (line 1) and reasonable composition (line 2), while Z-sampling achieves high prompt alignment.

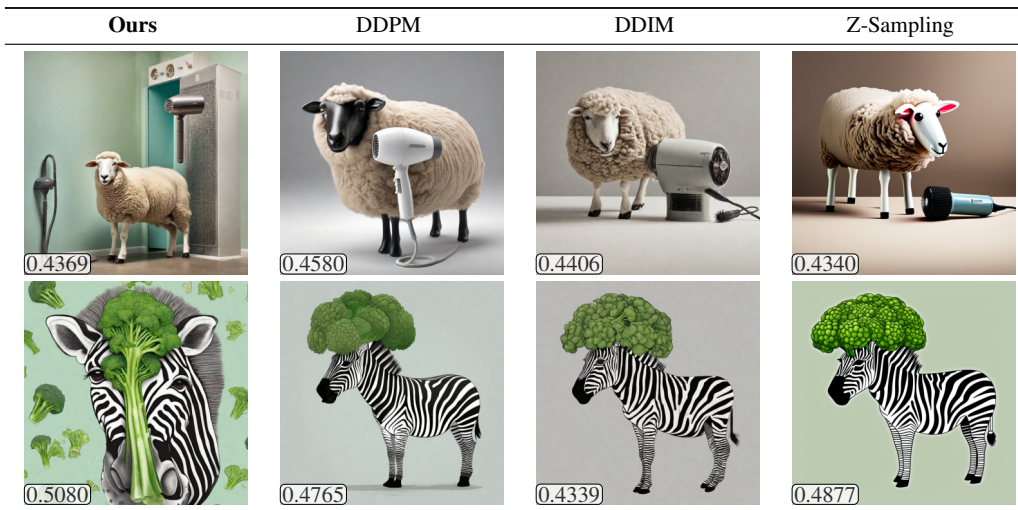
Q BROADER IMPACTS, LIMITATIONS AND FUTURE WORK

Q.1 BROADER IMPACTS

We outline the following potential broader impacts:

1. Our augmented RODE sampling provides an alternative solution for introducing randomness into ODE sampling without converting it into SDE sampling. This paradigm can po-

Figure 14: Failure cases of optimizing CLIP score.



tentially be adopted for RL with flow matching (Albergo & Vanden-Eijnden, 2022; Lipman et al., 2022; Liu et al., 2022) models;

- This study enhances the capabilities of pre-trained models in a training-free manner, such as enabling few-step inference performance to surpass that of multi-step inference in the same model. This suggests that pre-trained image generation models may not be fully utilized, prompting the community to explore the deeper potential of pre-trained models rather than solely focusing on training larger models.

Q.2 LIMITATIONS

We highlight the following limitations:

- The computational overhead of our method precludes its deployment in real-time scenarios;
- Our method attempt to fully utilized the pre-trained DMs, whose performance is bounded by the capabilities of the adopted base models;
- Our method is more effective in multi-step scenario, because the performance of our RODE sampling is highly dependent on a sufficient number of usable process noise. Despite our augmentation strategy partially alleviates this problem, we empirically observe dramatic performance degradation in very few-step inference scenarios;
- The *NFE dynamics* required for the simulation phase of MCTS increases linearly with the number of inference steps, leading to higher computational demands in scenarios with larger step counts;
- When using dense rewards from latent reward shaping, the extensive queries of intermediate rewards necessitate a large number of final reward calculation calls. This can become a notable efficiency bottleneck when reward function queries are costly;
- In the context of training-free DM alignment, our approach often lags behind gradient-based methods. This is because the later — which optimize the latent variables or noise prediction with gradient — can modify the latent representations to a greater extent and with higher accuracy, thereby holding greater promise for achieving superior performance.

Q.3 FUTURE WORK

We list the following potential directions for future work:

- Generalize to stochastic samplers other than DDIM (Song et al., 2020);

- 3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
2. Leverage advanced heuristic techniques to accelerate the search process for the optimal control sequence, thereby improving efficiency. For example, a problem-specific admissible heuristic that predicts the maximum attainable reward along a path can be introduced to prune unpromising branches at an early stage;
 3. Propose a global-local search strategy. Specifically, first conduct a global search that treating ϵ_{t_i} s as actions to identify high-reward regions. Then, perform a local search with η_{t_i} actions to fine-tune the denoising trajectory to release the potential of an initial noise to the maximum extent;
 4. Explore other parameter-efficient scaling methods;
 5. While we have used low-dimensional action sequences to easily visualize the high-dimensional denoising trajectories of DMs, the relationship between sampling results and action sequences requires further investigation.