

---

# ART for Diffusion Sampling: A Reinforcement Learning Approach to Timestep Schedule

---

Yilie Huang<sup>1</sup> Wenpin Tang<sup>1</sup> Xunyu Zhou<sup>1</sup>

## Abstract

We consider time discretization for score-based diffusion models to generate samples from a learned reverse-time dynamic on a finite grid. Uniform and hand-crafted grids can be suboptimal given a budget on the number of time steps. We introduce Adaptive Reparameterized Time (ART) that controls the clock speed of a reparameterized time variable, leading to a time change and uneven timesteps along the sampling trajectory while preserving the terminal time. The objective is to minimize the aggregate error arising from the discretized Euler scheme. We derive a randomized control companion, ART-RL, and formulate time change as a continuous-time reinforcement learning (RL) problem with Gaussian policies. We then prove that solving ART-RL recovers the optimal ART schedule, which in turn enables practical actor-critic updates to learn the latter in a data-driven way. Empirically, based on the official EDM pipeline, ART-RL improves Fréchet Inception Distance on CIFAR-10 over a wide range of budgets and transfers to AFHQv2, FFHQ, and ImageNet without the need of retraining.

## 1. Introduction

Diffusion models (Ho et al., 2020; Song & Ermon, 2019; Song et al., 2021b) are a family of generative models that create samples from unknown target distributions based on given samples. They have played a key role in the success in text-to-image creators such as DALL-E 2 (Ramesh et al., 2022) and Stable Diffusion (Rombach et al., 2022), in text-to-video generators such as Sora (OpenAI, 2024), Make-A-Video (Singer et al., 2023) and Veo (Google, 2024) and, more recently, in diffusion large language models such as Mercury (Khanna et al., 2025) and LLaDA (Nie et al., 2025).

<sup>1</sup>Department of Industrial Engineering and Operations Research, Columbia University, New York, NY, USA. Correspondence to: Yilie Huang <yh2971@columbia.edu>.

A diffusion model consists of two steps: learn the target score function by training on given samples (pretraining), and generate new samples from the learned model (inference). The focus of this paper is on diffusion sampling/inference, where it involves choosing timesteps or time discretization of the learned model. Existing work uses either uniform or hand-crafted timesteps (Song et al., 2021a;b; Karras et al., 2022; Chen et al., 2023a), lacking a principled framework to schedule timesteps. Here our purpose is to provide a control-theoretic framework, which allows strategic and systematic time discretization in diffusion sampling in order to minimize the aggregate error when applying the Euler scheme. The main contribution of this paper is as follows:

- *Methodology*: We propose an optimal control framework – *Adaptive Reparameterized Time* (ART), to schedule timesteps for diffusion sampling. The idea is to treat the speed of the diffusion sampler as control to reparameterize time and adaptively redistribute computation along the sampling trajectory. To solve the ART problem, we propose a continuous-time reinforcement learning (CTRL) approach – ART-RL, premised upon the recent development in Wang et al. (2020); Jia & Zhou (2022a;b).
- *Theory*: We provide a rigorous theoretical analysis of the CTRL formulation in the ART setting. We show that the optimal ART control coincides with the mean of the optimal randomized Gaussian policy in ART-RL, establishing a precise link between ART and ART-RL. Based on this connection, we derive actor-critic update rules for learning adaptive time schedules.
- *Experiments*: Our ART-RL schedule consistently outperforms uniform and EDM schedules (Karras et al., 2022) within the official EDM pipeline. It improves Fréchet Inception Distance (FID) across a broad range of sampling time budgets on CIFAR-10, with particularly strong gains at low budgets. Moreover, under the scheduled timesteps for CIFAR-10, the improvement transfers without retraining to other datasets such as AFHQv2, FFHQ, and ImageNet.

a principled approach to scheduling timesteps for generative diffusion sampling. The proposed ART-RL method is purely data-driven, and shows robust performance and generalizations empirically.

**Relevant literature:** Diffusion models as generative tools were first proposed by Ho et al. (2020) (DDPM) and Song et al. (2021a) (DDIM) in the discrete setting. The pioneering work of Song et al. (2021b) introduces a continuous-time formulation of diffusion models, providing a unified treatment that encompasses earlier discrete-time models. Various sampling methods have been proposed for diffusion inference, including predictor-corrector sampler (Song et al., 2021b), exponential integrator (Zhang & Chen, 2023), and higher-order solvers (Zhang et al., 2023; Wu et al., 2024). See also Lee et al. (2022); Chen et al. (2023a;b); Li et al. (2024); Benton et al. (2024); Li & Yan (2024); Huang et al. (2025a) for convergence analysis of diffusion models (with either uniform or hand-crafted timesteps).

Continuous-time reinforcement learning (CTRL) was first formulated by Wang et al. (2020), where exploration is modeled via stochastic relaxed control to capture the trial-and-error nature of reinforcement learning. Subsequent work developed a model-free theoretical foundation via martingale methods (Jia & Zhou, 2022a;b; 2023; Tang & Zhou, 2024), established performance guarantees (Huang et al., 2024; 2025b), and studied policy optimization (Zhao et al., 2023). The CTRL framework has also been applied to training and fine-tuning diffusion models in generative AI (Gao et al., 2024; Zhao et al., 2024; 2025).

**Organization:** Section 2 reviews diffusion models and introduces the ART control formulation. Section 3 describes how CTRL is used to solve ART, leading to the ART-RL approach. Section 4 presents the ART-RL algorithm, followed by empirical results in Section 5. Section 6 concludes. The proofs and additional numerical results are placed in the appendix.

## 2. Diffusion Models and ART

### 2.1. Continuous-Time Score-Based Diffusion Models

We briefly review continuous-time score-based diffusion models; see Tang & Zhao (2025) for a recent survey. Given samples from an unknown target distribution, a forward process corrupts data over time  $\tau \in [0, T]$  toward a tractable reference distribution, and sampling proceeds by integrating a learned reverse-time process.

**Diffusion sampling.** A forward process obeys the Itô SDE

$$d\bar{x}(\tau) = -f(\tau)\bar{x}(\tau)d\tau + g(\tau)dw(\tau), \quad (1)$$

where  $\tau \in [0, T]$ ,  $\bar{x}(0) \sim p_0$  is the target data distribution,  $w$  is a standard Wiener process in  $\mathbb{R}^d$ , and  $f$  and  $g$  are coefficient functions. Let  $p_\tau$  be the law of  $\bar{x}(\tau)$  and  $S(\tau, x) = \nabla_x \log p_\tau(x)$  be the score.

For sampling, we use the probability flow ODE associated with (1), which shares the same marginals as the reverse-time SDE; see e.g. Tang & Zhao (2025, Theorem 5.1). Denote the backward state by  $\tilde{x}(\tau) := \bar{x}(T - \tau)$  with initialization  $\tilde{x}(0) \sim p_T$ . Replacing  $S$  with a trained score model  $\hat{S}$  yields

$$\frac{d\tilde{x}(\tau)}{d\tau} = f(T - \tau)\tilde{x}(\tau) + \frac{1}{2}g(T - \tau)^2\hat{S}(T - \tau, \tilde{x}(\tau)), \quad (2)$$

where  $\tau \in [0, T]$ .

**Euler discretization.** We integrate (2) on a grid  $0 = \tau_0 < \tau_1 < \dots < \tau_K = T$  with step sizes  $h_i = \tau_{i+1} - \tau_i$ , and write  $\tilde{x}_i := \tilde{x}(\tau_i)$ . The explicit Euler update is

$$\tilde{x}_{i+1} = \tilde{x}_i + h_i \left[ f(T - \tau_i)\tilde{x}_i + \frac{1}{2}g(T - \tau_i)^2\hat{S}(T - \tau_i, \tilde{x}_i) \right]. \quad (3)$$

A uniform grid  $\tau_i = iT/K$  is simple but allocates evaluations uniformly even when the reverse dynamics vary along the trajectory. A simple intuition is that early stages, where samples resemble noise, may tolerate coarser steps, while later stages will benefit from finer resolutions. Such considerations motivate adaptive, data-driven schedules that redistribute steps under a given, fixed total time budget  $T$ .

### 2.2. ART: Time Reparameterization via Control

Motivated by the drawback of uniform timesteps in backward sampling discussed above, we introduce a reparameterized sampling clock and formalize adaptive timestep selection in the reverse process as a control problem. The key idea is to replace the uniform progression of original time by an adaptive, controlled progression that can accelerate in some segments and decelerate in others, thereby redistributing computational effort along the trajectory.

To this end, let  $\psi : [0, T] \rightarrow \mathbb{R}$  be a continuous time mapping from the new clock  $t$  to the original diffusion time  $\tau$  (i.e.,  $\tau = \psi(t)$ ), with  $\psi(0) = 0$  and  $\psi(T) = T$ . On this new clock we write the state as  $x(t) := \tilde{x}(\psi(t))$ , namely the backward state evaluated at diffusion time  $\psi(t)$ , with  $x(0) \sim p_T$ . We define the control as  $\theta(t) := \dot{\psi}(t)$ , which quantifies the instantaneous rate of change in diffusion time with respect to  $t$  and satisfies  $\int_0^T \theta(t) dt = \psi(T) - \psi(0) = T$ . To keep the formulation general, we do not impose monotonicity of  $\psi$ , so  $\theta$  may take either sign. The monotone time-change case (equivalently,  $\theta(t) \geq 0$  almost everywhere) is included as a special case of this formulation. In particular, if we discretize the new clock uniformly as

$0 = t_0 < t_1 < \dots < t_K = T$ , then the induced original-time grid is  $\tau_i := \psi(t_i)$  with nonuniform step sizes  $\tau_{i+1} - \tau_i$ . Intuitively,  $x(\cdot)$  tracks the generative trajectory under the reparameterized time, while  $\theta(\cdot)$  determines how quickly the trajectory advances in original time per unit step on the new clock (larger  $\theta(t)$  corresponds to faster progression at  $t$ ). We refer to this time reparameterized sampling framework as *Adaptive Reparameterized Time* (ART).

We next state the controlled dynamics induced by ART. By the chain rule, with  $x(t) = \tilde{x}(\psi(t))$ , we have

$$\begin{cases} \dot{x}(t) = \theta(t)F(x(t), \psi(t)), & x(0) \sim p_T, \\ \dot{\psi}(t) = \theta(t), & \psi(0) = 0, \psi(T) = T, \end{cases} \quad (4a)$$

$$(4b)$$

where the backward probability-flow vector evaluated at original time  $T - \psi$  is

$$F(x, \psi) := f(T - \psi)x + \frac{1}{2}g(T - \psi)^2 \hat{S}(T - \psi, x). \quad (5)$$

Recall the control  $\theta(t) = \dot{\psi}(t)$  is the local time-scaling factor (time-warping rate) with respect to the new clock and satisfies the total-time constraint

$$\int_0^T \theta(t) dt = \psi(T) - \psi(0) = T, \quad (6)$$

which represents the overall amount of time progression to be allocated along the trajectory.

To assess how Euler behaves on the  $t$ -clock, we relate its local approximation error over a step to the curvature of the right-hand side in (4a). We proceed analogously to the Euler discretization in (3). On an interval  $t \in [t_i, t_{i+1})$  with stepsize  $h_i := t_{i+1} - t_i$ , the Euler update uses one control value per step; we denote this step value by  $\theta_i$ .

Define the one-step Euler error proxy on the  $t$ -clock by

$$E_i := x(t_{i+1}) - \left( x(t_i) + h_i \theta_i F(x(t_i), \psi(t_i)) \right).$$

A Taylor expansion yields

$$E_i = \frac{h_i^2}{2} \theta_i^2 Q(x(t_i), \psi(t_i)) + O(h_i^3), \quad (7)$$

where

$$\begin{aligned} Q(x, \psi) &= A(s, x) B(s, x) - g(s)g'(s)\hat{S}(s, x) \\ &\quad - f'(s)x - \frac{1}{2}g(s)^2 \partial_s \hat{S}(s, x), \end{aligned} \quad (8)$$

with  $s := T - \psi$ ,  $A(s, x) := f(s)I_d + \frac{1}{2}g(s)^2 \nabla_x \hat{S}(s, x)$ , and  $B(s, x) := f(s)x + \frac{1}{2}g(s)^2 \hat{S}(s, x)$ . Thus, to leading order, the magnitude of the Euler local error on step  $i$  is proportional to  $\theta_i^2 |Q(x(t_i), \psi(t_i))|$ , where  $|\cdot|$  denotes the Euclidean norm throughout. In implementation,  $\nabla_x \hat{S}(s, x)$  is

never formed explicitly; it is only queried through Jacobian-vector product computed by automatic differentiation.

This motivates using  $|Q(x, \psi)| \theta(t)^2$  as an error-density surrogate for allocating the time-warping rate. With a fixed time budget, we choose  $\theta = \theta(\cdot)$  to minimize the overall residual surrogate subject to the constraint  $\int_0^T \theta(t) dt = T$ , enforced via a Lagrange multiplier  $\gamma \in \mathbb{R}$ . The resulting objective is to maximize

$$\begin{aligned} J^\theta(s, y, \phi) &= \mathbb{E} \left[ \int_s^T (-|Q(x(t), \psi(t))| \theta^2(t) \right. \\ &\quad \left. - \gamma \theta(t)) dt + \gamma T \mid x(s) = y, \psi(s) = \phi \right]. \end{aligned} \quad (9)$$

Denote the optimal value function to be

$$V(s, y, \phi) = \max_{\theta=\theta(\cdot)} J^\theta(s, y, \phi). \quad (10)$$

In summary, ART recasts timestep allocation as continuous-time control of the time-warping rate  $\theta$  under (4) and (9). The next section develops an RL-based procedure (ART-RL) to learn  $\theta$  in a data-driven way.

### 3. Randomized Control and RL Formulation

The ART problem just formulated has no closed-form solution, not is its Hamilton-Jacobi-Bellman (HJB) equation on a high-dimensional state space in  $x \in \mathbb{R}^d$  numerically intractable due to the curse of dimensionality. We instead resort to an RL-based solution by considering a randomized version of the problem in which  $\theta$  is generated by a stochastic policy. Unlike most RL work, the randomization here is not for exploration; rather it is a *technical* device that embeds the underlying optimal control problem into the continuous-time RL framework that has been well developed recently in both theory and algorithms. We call this approach *Adaptive Reparameterized Time via Reinforcement Learning* (ART-RL), and present it in this section.

#### 3.1. ART-RL: Auxiliary Problem with Gaussian Policies

We model control randomization with a stochastic policy that, for each triple  $(t, x, \psi)$ , specifies a probability distribution of the time warping rate  $\theta$ . Inspired by the entropy-regularized formulation in the discrete-time setting (Ziebart et al., 2008) and its continuous-time counterpart (Huang et al., 2022), where Gaussian policies are optimal, we consider the Gaussian family indexed by  $\lambda \geq 0$ :

$$\pi^{(\lambda)}(\cdot \mid t, x, \psi) = \mathcal{N}\left(\mu(t, x, \psi), \frac{\lambda}{|Q(x, \psi)|}\right), \quad (11)$$

where  $\mu$  is a measurable deterministic policy and  $Q$  is defined in (8). Let  $\Pi^{(\lambda)}$  denote the class of policies of the form

(11) with finite second moments. This parameterization adapts the variance to the problem geometry: since  $Q(x, \psi)$ , as shown in (7), is proportional to the local truncation error of the Euler-discretized probability-flow dynamics, scaling the variance inversely with  $|Q|$  reduces randomness in stiff regions while allowing more randomness elsewhere. The scalar  $\lambda$  controls the overall noise level without affecting the mean. For analysis we assume  $|Q(x, \psi)| > 0$  almost surely on compact intervals; in practice we replace  $|Q|$  by  $|Q| \vee \varepsilon$  for a small  $\varepsilon > 0$ .

Given a policy  $\pi^{(\lambda)} \in \Pi^{(\lambda)}$ , let  $(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))_{t \in [0, T]}$  denote the corresponding state process which, according to Wang et al. (2020), satisfy

$$\begin{cases} \dot{x}^{\pi^{(\lambda)}}(t) = \int_{\mathbb{R}} \theta F(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)) \pi_t^{(\lambda)}(\theta) d\theta & (12a) \\ x^{\pi^{(\lambda)}}(0) = x_0 \sim P_T, \\ \dot{\psi}^{\pi^{(\lambda)}}(t) = \int_{\mathbb{R}} \theta \pi_t^{(\lambda)}(\theta) d\theta, & (12b) \\ \psi^{\pi^{(\lambda)}}(0) = 0, \psi^{\pi^{(\lambda)}}(T) = T. \end{cases}$$

The corresponding total-time constraint is  $\int_0^T \int_{\mathbb{R}} \theta \pi_t^{(\lambda)}(\theta) d\theta dt = T$ . The objective function for the auxiliary problem is

$$\begin{aligned} J^{\pi^{(\lambda)}}(s, y, \phi) = \mathbb{E} \Big[ & \gamma T + \lambda T + \int_s^T \int_{\mathbb{R}} \pi_t^{(\lambda)}(\theta) \\ & (-|Q(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))|^2 - \gamma \theta) d\theta & (13) \\ & dt \mid x^{\pi^{(\lambda)}}(s) = y, \psi^{\pi^{(\lambda)}}(s) = \phi \Big], \end{aligned}$$

where the Lagrange multiplier  $\gamma \in \mathbb{R}$  enforces the total-time constraint and  $\lambda \geq 0$  is the variance parameter from (11). The optimal value function is

$$V^{(\lambda)}(s, y, \phi) = \max_{\pi^{(\lambda)} \in \Pi^{(\lambda)}} J^{\pi^{(\lambda)}}(s, y, \phi). \quad (14)$$

### 3.2. Relationship between ART and ART-RL

We now establish the relationship between the ART control problem (10) and the ART-RL auxiliary problem (14). By dynamic programming, the optimal value function  $V$  for (10) satisfies the HJB equation

$$V_t + \sup_{\theta} \left\{ (V_x^\top F(x, \psi) + V_\psi - \gamma) \theta - |Q(x, \psi)| \theta^2 \right\} = 0, \quad (15)$$

with terminal condition  $V(T, x, \psi) = \gamma T$ .

Moreover, by Wang et al. (2020), the optimal value function

$V^{(\lambda)}$  of (14) satisfies the exploratory HJB equation:

$$\begin{aligned} V_t^{(\lambda)} + \sup_{\mu} \Big\{ & (V_x^{(\lambda)\top} F(x, \psi) + V_\psi^{(\lambda)} - \gamma) \mu \\ & - |Q(x, \psi)| \left( \mu^2 + \frac{\lambda}{|Q(x, \psi)|} \right) \Big\} = 0 \end{aligned} \quad (16)$$

with terminal condition  $V^{(\lambda)}(T, x, \psi) = (\gamma + \lambda)T$ .

The connection between ART and ART-RL is characterized by the following two theorems.

**Theorem 3.1** (Value function shift). *If  $V$  is a classical solution of (15), then*

$$V^{(\lambda)}(t, x, \psi) = V(t, x, \psi) + \lambda t \quad (17)$$

*is a classical solution of (16).*

**Theorem 3.2** (Recovery of the optimal ART control). *Under the conditions of Theorem 3.1, define*

$$\mu^*(t, x, \psi) = \frac{V_x^\top F(x, \psi) + V_\psi - \gamma}{2|Q(x, \psi)|}. \quad (18)$$

*Then the Gaussian policy*

$$\pi^{(\lambda)*}(\cdot \mid t, x, \psi) = \mathcal{N}\left(\mu^*(t, x, \psi), \frac{\lambda}{|Q(x, \psi)|}\right) \quad (19)$$

*is an optimal policy for the auxiliary problem (14) subject to the dynamics (12). Furthermore,  $\mu^*$  is an optimal policy for the original problem (10) subject to the dynamics (4).*

Proofs of Theorems 3.1 and 3.2 are provided in Appendix A.

Theorems 3.1 and 3.2 show that the ART solution can be retrieved by solving the ART-RL problem (14). The latter can be solved using an actor-critic scheme that does not work directly on the original control problem (10). Specifically, the Gaussian policy class makes policy evaluation and improvement tractable; details are given in the next section.

## 4. ART-RL Actor-Critic Algorithm

By Theorems 3.1 and 3.2, it suffices to solve the auxiliary ART-RL problem in order to obtain the original ART solution.

### 4.1. Algorithm Design

Our algorithm builds on the continuous-time actor-critic framework of Jia & Zhou (2022b), adapted to ART-RL with randomized Gaussian control. We use two neural networks  $NN^{\vartheta_c}$  and  $NN^{\vartheta_a}$  for the value function and policy mean, and parameterize the value function and policy as

$$\begin{aligned} \hat{V}^{\vartheta_c}(t, x, \psi) &= NN^{\vartheta_c}(t, x, \psi) + \lambda t, \\ \hat{\pi}^{\vartheta_a}(\cdot \mid t, x, \psi) &= \mathcal{N}\left(NN^{\vartheta_a}(t, x, \psi), \frac{\lambda}{|Q(x, \psi)|}\right). \end{aligned} \quad (20)$$



Let  $\theta^{\hat{\pi}^{\vartheta_a}} \sim \hat{\pi}^{\vartheta_a}$  be a Gaussian control process with corresponding (observable) state process  $(x^{\theta^{\hat{\pi}^{\vartheta_a}}}(t), \psi^{\theta^{\hat{\pi}^{\vartheta_a}}}(t))$ . For brevity, in (21) we write  $(x(t), \psi(t), \theta(t)) = (x^{\theta^{\hat{\pi}^{\vartheta_a}}}(t), \psi^{\theta^{\hat{\pi}^{\vartheta_a}}}(t), \theta^{\hat{\pi}^{\vartheta_a}}(t))$ . By Jia & Zhou (2022b),  $\hat{V}^{\vartheta_c}$  and  $\hat{\pi}^{\vartheta_a}$  satisfy the coupled moment conditions

$$\begin{aligned} \mathbb{E} \left[ \int_0^T \frac{\partial N N^{\vartheta_c}(t, x(t), \psi(t))}{\partial \vartheta_c} \left( d\hat{V}^{\vartheta_c}(t, x(t), \psi(t)) \right. \right. \\ \left. \left. - (|Q(x(t), \psi(t))| \theta(t)^2 + \gamma \theta(t)) dt \right) \right] = 0, \\ \mathbb{E} \left[ \int_0^T \left( \frac{\partial \log \hat{\pi}^{\vartheta_a}(\theta(t) | t, x, \psi)}{\partial \vartheta_a} \right) \left( d\hat{V}^{\vartheta_c}(t, x(t), \right. \right. \\ \left. \left. \psi(t)) - (|Q(x(t), \psi(t))| \theta(t)^2 + \gamma \theta(t)) dt \right) \right] = 0. \end{aligned} \quad (21)$$

This leads to standard moment conditions in the RL literature (Sutton & Barto, 1998; Huang & Zhou, 2025). The learnable parameters  $\vartheta_c$  and  $\vartheta_a$  can then be updated via stochastic approximation, as shown in the next subsection.

#### 4.2. Implementable ART-RL Algorithm

We index iterations by  $n$ ; for instance,  $\vartheta_{c,n}$  denotes the iterate for  $\vartheta_c$  at iteration  $n$ . At each iteration, a trajectory  $(x_n, \psi_n, \theta_n)$  is generated under the current policy  $\hat{\pi}^{\vartheta_{a,n}}$ , and parameters are updated with learning rate  $a_n > 0$ . In implementation, time is discretized on a uniform grid  $0 = t_0 < t_1 < \dots < t_K = T$  with  $\Delta t = T/K$ . For the  $n$ -th iteration, write  $\hat{V}_k^{\vartheta_{c,n}} := \hat{V}^{\vartheta_{c,n}}(t_k, x_n(t_k), \psi_n(t_k))$ . Viewing the moment conditions (21) as equations in  $(\vartheta_c, \vartheta_a)$ , we apply stochastic approximation and Riemann discretization to obtain implementable critic and actor updates.

$$\begin{aligned} \vartheta_{c,n+1} \leftarrow \vartheta_{c,n} + a_n \sum_{k=0}^{K-1} \frac{\partial N N^{\vartheta_{c,n}}}{\partial \vartheta_c}(t_k, x_n(t_k), \psi_n(t_k)) \times \\ \left[ \hat{V}_{k+1}^{\vartheta_{c,n}} - \hat{V}_k^{\vartheta_{c,n}} - \gamma_n \theta_n(t_k) \Delta t \right. \\ \left. - |Q(x_n(t_k), \psi_n(t_k))| \theta_n(t_k)^2 \Delta t \right], \end{aligned} \quad (22a)$$

$$\begin{aligned} \vartheta_{a,n+1} \leftarrow \vartheta_{a,n} + a_n \sum_{k=0}^{K-1} \frac{\partial \log \hat{\pi}_{t_k}^{\vartheta_{a,n}}(\theta_n(t_k))}{\partial \vartheta_a} \times \\ \left[ \hat{V}_{k+1}^{\vartheta_{c,n}} - \hat{V}_k^{\vartheta_{c,n}} - \gamma_n \theta_n(t_k) \Delta t \right. \\ \left. - |Q(x_n(t_k), \psi_n(t_k))| \theta_n(t_k)^2 \Delta t \right], \end{aligned} \quad (22b)$$

where we write  $\hat{\pi}_{t_k}^{\vartheta_{a,n}}(\theta) := \hat{\pi}^{\vartheta_{a,n}}(\theta | t_k, x_n(t_k), \psi_n(t_k))$ . The update for the Lagrange multiplier is

$$\gamma_{n+1} \leftarrow \gamma_n + a_n (\psi_n(T) - T), \quad (23)$$

where  $\psi_n(T) = \psi_n(t_K)$  is the terminal state on the grid. We summarize the resulting time-discretized ART-RL actor-

critic scheme in Algorithm 1. The inner loop generates one trajectory under the current Gaussian policy (20), and the outer loop then updates the actor, critic, and Lagrange multiplier via (22) and (23).

---

#### Algorithm 1 Time-discretized ART-RL Actor-Critic

---

```

for  $n = 1$  to  $N$  do
  Set  $k = 0, t = t_k = 0$ , initialize  $(x_n(t_0), \psi_n(t_0))$ 
  while  $t < T$  do
    Sample control  $\theta_n(t_k)$  from the Gaussian policy (20)
    Update the states by the ART dynamics (4)
    Increment time:  $t_{k+1} = t_k + \Delta t, k \leftarrow k + 1$ 
  end while
  Collect trajectory  $\{(t_k, x_n(t_k), \psi_n(t_k), \theta_n(t_k))\}_{k=0}^{K-1}$ 
  Update critic and actor parameters via (22a) and (22b)
  Update the Lagrange multiplier  $\gamma_{n+1}$  via (23)
end for

```

---

## 5. Numerical Experiments

We now evaluate ART-RL across experiments varying in dimensionality, model capacity, and experimental protocol. Code is provided in the supplementary material.

**Datasets.** We consider both synthetic and real-image settings. In a one-dimensional experiment, we construct a synthetic target distribution on  $\mathbb{R}$  with a known score function, isolating the effect of time reparameterization from score-estimation error. For high-dimensional image generation, we use CIFAR-10 (Krizhevsky & Hinton, 2009), AFHQv2 (Choi et al., 2020), FFHQ (Karras et al., 2019), and ImageNet (Russakovsky et al., 2015) under the EDM pipeline (Karras et al., 2022).

**Timestep schedules and baselines.** We compare three timestep schemes for all the experiments. The first is *Uniform*, which uses an equally spaced grid in the physical time variable  $\tau \in [0, T]$  and serves as a simple baseline.

The second is *EDM*, the hand-crafted schedule of Karras et al. (2022), which is widely used in diffusion models and performs strongly on standard image benchmarks. In our implementation, the discrete timesteps are given by

$$\tau_k = \left( \sigma_{\max}^{1/\rho} + \frac{k}{K} (\sigma_{\min}^{1/\rho} - \sigma_{\max}^{1/\rho}) \right)^\rho, \quad k = 0, \dots, K,$$

with  $\sigma_{\min} > 0, \sigma_{\max} > \sigma_{\min}$ , and exponent  $\rho > 0$  (following Karras et al. 2022, we set default  $\rho = 7$ ). This can be viewed as placing a uniform grid in the transformed coordinate  $\sigma^{1/\rho}$  between  $\sigma_{\max}^{1/\rho}$  and  $\sigma_{\min}^{1/\rho}$ , which corresponds to a fixed, hand-designed time reparameterization.

The third scheme, *ART-RL*, is our learned schedule produced by Algorithm 1. In the ART formulation, the control  $\theta$

induces a time change  $\psi$ , and we place a uniform grid in the reparameterized  $t$ -clock. When  $\psi(t) = t$ , the grid reduces to the Uniform scheme. When  $\psi(t)$  is proportional to  $\sigma^{1/\rho}$ , with  $\sigma$  following the EDM noise schedule, the induced grid coincides with the EDM schedule up to a constant rescaling. By contrast, ART or ART-RL optimizes  $\psi$  to minimize Euler discretization error on the reparameterized clock, and can therefore outperform these hand-crafted schemes by adapting to other geometries when beneficial.

**Evaluation metrics.** In the one-dimensional setting, we evaluate backward sampling quality using the squared Wasserstein distance ( $W_2$ ) between the empirical distribution of generated samples and the target distribution, together with the number of timesteps used by the Euler discretization of the probability flow ODE (no neural network is involved). For the image experiments, we follow standard practice and report the Fréchet Inception Distance (FID) as a function of the number of function evaluations (NFEs). In the EDM pipeline (Karras et al., 2022), the compared methods differ only in timestep schedule; so the FID–NFE curves isolate the effect of time grids.

### 5.1. One-Dimensional Example with Analytical Score

We consider a one-dimensional example where the score function is available in closed form; so the only source of approximation error is the timestep schedule. The forward process starts from  $p_0 = \mathcal{N}(0, 1)$  and follows the Itô SDE  $dx(t) = \sqrt{2t} dw(t)$  over the horizon  $T = 3$ . In this case  $x(t) \sim \mathcal{N}(0, 1 + t^2)$ , yielding the terminal distribution  $p_T = \mathcal{N}(0, 10)$  and the *exact* score  $S(t, x) = -x/(1 + t^2)$ .

Substituting this score into (5) and (8) gives  $F(x, \psi) = -(T - \psi)x/(1 + (T - \psi)^2)$  and  $Q(x, \psi) = x/(1 + (T - \psi)^2)^2$ . This setup isolates the effect of timestep schedules without score-estimation error.

We next examine the learned control  $\theta$ . During training of the ART-RL Algorithm 1 with  $K = 100$  timesteps, we record the executed  $\theta$ -values along the last 10,000 backward trajectories. For each trajectory, we normalize the realized  $\theta$ -sequence so that the induced total time change sums to  $T$ , removing any residual over- or under-shoot of  $\psi(T)$ . From these normalized trajectories, we compute at each timestep the empirical mean of  $\theta$  and the interquartile range (IQR, 25–75 percentiles), shown in Figure 1.

Figure 1 shows that the mean curve of  $\theta$  is smooth with an extremely narrow IQR band. Moreover, the 99% confidence band (Appendix B.1, Figure 3) is visually indistinguishable from the mean. This suggests that, in this example, the learned control  $\theta$  depends only weakly on the state and can be effectively regarded as a function of time alone; that is, the policy collapses to an almost time-only schedule.

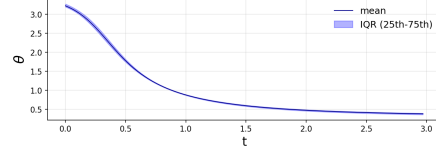


Figure 1. Empirical mean (solid line) and 25–75 percent interquartile range (shaded region) of the learned control  $\theta$  across time.

Motivated by this observation, we perform a simple distillation step in our sampling procedure by discarding the neural-network actor and replacing it with the empirical mean curve of  $\theta$  as a fixed function of  $t$ . This distilled schedule has two advantages.

First, distillation removes the cost of per-step policy computation, including evaluating the actor network and the state-dependent variance via  $Q$ . Although these components are much cheaper than the score model, repeated evaluation along sampling trajectories still incurs nontrivial overhead. After distillation, ART-RL sampling requires no extra computation beyond standard schemes such as Uniform or EDM; the timestep sequence is precomputed and reused.

Second, it eliminates residual mismatch in the terminal time. While the learned actor enforces  $\psi(T) \approx T$  in expectation, individual trajectories may slightly overshoot or undershoot  $T$  when  $\theta$  is produced by a neural network at every timestep. This discrepancy is negligible for small  $K$  but becomes more pronounced as  $K$  grows and finer time resolution matters. By distilling to a deterministic schedule whose increments are normalized to sum to  $T$ , we guarantee that the induced grid hits  $T$  exactly, improving the numerical fidelity of the discretized probability flow ODE.

In view of these advantages, we use the distilled ART-RL schedule for sampling in all our experiments, while the schedule is still trained by Algorithm 1. We now report the numerical results for the one-dimensional example.

Table 1.  $W_2$  vs. timesteps  $K$  in the one-dimensional experiment.

$K$	2	5	10	20	50	100
Uniform	.468	.215	.114	.060	.027	.016
EDM-T	.414	.195	.105	.056	.025	.015
EDM	.664	.319	.177	.094	.041	.023
ART-RL	<b>.345</b>	<b>.149</b>	<b>.079</b>	<b>.042</b>	<b>.020</b>	<b>.013</b>

The EDM paper (Karras et al., 2022) recommends the default exponent  $\rho = 7$  and reports strong performance on image datasets under their sampling pipeline. However, in this one-dimensional setting the same choice performs poorly: Table 1 shows that EDM with  $\rho = 7$  is significantly worse than all the other methods across all timestep counts and even underperforms the Uniform grid. This highlights

an important point: although EDM provides a training-free schedule, its effectiveness is *problem dependent*, and the choice of hyperparameters can matter substantially.

For comparison, we also report an EDM variant with a nondefault choice of  $\rho = 3$  that yields clearly improved performance after tuning, which we label “EDM-T”. The tuned EDM schedule consistently outperforms Uniform, confirming that EDM can improve sampling quality when appropriately calibrated. But ART-RL achieves consistently lower W2 errors than both Uniform and EDM-T, with a substantial margin across all  $K$ . This suggests that while EDM can be fine-tuned, such tuning incurs additional cost and still does not match the accuracy of ART-RL, whose schedule is learned directly from the underlying dynamics.

## 5.2. CIFAR-10 under the EDM Pipeline

We use the official EDM pipeline (Karras et al., 2022) and keep all details fixed (score model, solver, noise conditioning, and hyperparameters). ART-RL replaces only the time grid. Since our objective is theoretically motivated by an Euler local error proxy, we additionally run a solver-aligned Euler ablation as a controlled check. ART-RL improves FID over both EDM and Uniform across all tested budgets; for example, at NFE = 5 it attains 28.16 versus 49.10 (EDM) and 214.60 (Uniform), and at NFE = 30 it attains 4.06 versus 4.21 (EDM) and 85.83 (Uniform); see Appendix B.2 for the full sweep.

We report the main image results with the default Heun solver to follow the standard EDM setup and to demonstrate that the learned time grid is not specific to Euler updates and remains effective under higher-order solvers. Precisely, we test  $K \in \{2, 5, 7, 10, 18\}$  steps and use Euler only for the final step update, matching the EDM implementation and giving NFE =  $2K - 1$ . EDM reports its best result at  $K = 18$  (NFE = 35); so we report results up to  $K = 18$  for a budget-matched comparison.

Table 2. FID vs. NFE on CIFAR-10 under the EDM pipeline.

NFE	3	9	13	19	35
Uniform	280.29	213.13	191.69	168.87	118.02
EDM	465.83	35.54	6.79	2.54	<b>1.85</b>
ART-RL	<b>152.86</b>	<b>32.13</b>	<b>5.44</b>	<b>2.45</b>	<b>1.85</b>

The results in Table 2 show several noteworthy patterns. First, the default EDM schedule is not uniformly strong across budgets: at small NFE = 3, EDM performs poorly and even worse than the simple Uniform grid, whereas ART-RL already yields substantial improvement. From NFE = 9 onward, EDM starts to outperform Uniform significantly and enters the high-quality regimes.

Second, ART-RL consistently attains the lowest FID among the three schedules at all reported NFEs. The improvement is most pronounced at small and moderate budgets (NFE = 3, 9, 13), where ART-RL corrects the failure mode of EDM at very small step counts and further improves once EDM becomes competitive. At larger budgets (NFE = 19, 35), the gap between ART-RL and EDM narrows, with both achieving strong FID values and ART-RL never underperforming EDM.

Third, the visuals in Appendix B.3, Figure 4, corroborate the numerical results. The Uniform grid produces noticeably blurry images even at NFE = 35, whereas both EDM and ART-RL yield sharp samples at moderate and large budgets. At very small budgets (e.g., NFE = 3), ART-RL already produces recognizable objects, while EDM outputs remain close to noise, mirroring the quantitative gaps in Table 2.

In sum, these results indicate that our learned time parameterization can be used as a drop-in replacement inside a strong sampling pipeline such as EDM, substantially improving robustness and sample quality in low- and mid-computation regimes while maintaining superior performance when more function evaluations are available.

## 5.3. Generalization of the ART-RL Time Schedule

We now examine how well the ART-RL time schedule transfers beyond the configuration in which it was learned, both across different timestep counts and across datasets without retraining. In the following, we focus on EDM and ART-RL; the Uniform grid performs substantially worse in these settings and is omitted for clarity.

### 5.3.1. ROBUSTNESS ACROSS TIME GRIDS VIA INTERPOLATION AND EXTRAPOLATION

To investigate whether ART-RL time schedules can be robustly reused across different timestep counts, we use CIFAR-10 as a testbed. Both methods under comparison follow the same experimental setup as in Section 5.2. For ART-RL, we take the learned time schedule obtained under  $K = 18$  in Section 5.2 and construct new grids for  $K = 4, 6, 9, 12, 20$  by log-linear interpolation and extrapolation. For EDM, the timestep sequence at each  $K$  is computed directly from its analytic formula.

Table 3. FID vs. NFE on CIFAR-10 for interpolated and extrapolated timestep counts.

NFE	7	11	17	23	39
EDM	85.80	14.42	3.11	2.06	<b>1.85</b>
ART-RL	<b>33.73</b>	<b>6.59</b>	<b>2.57</b>	<b>2.00</b>	<b>1.85</b>

Table 3 shows that the ART-RL schedule generalizes smoothly to both interpolated and extrapolated timestep

counts. Across all NFEs, ART-RL preserves the performance hierarchy observed in Section 5.2: its FID is substantially lower than EDM at small and moderate NFEs and remains at least as good at higher budgets, suggesting that the schedule learned at  $\text{NFE} = 35$  captures a stable structure that persists under these modifications. Corresponding images are shown in Appendix B.4, Figure 5.

### 5.3.2. CROSS-DATASET TRANSFER: AFHQv2, FFHQ, AND IMAGENET

We next investigate whether the ART-RL time schedule learned on CIFAR-10 can transfer directly to other datasets without retraining. For AFHQv2, FFHQ, and ImageNet, we keep all implementation details in the EDM pipeline unchanged, and replace only the timestep grid with the ART-RL schedule trained from CIFAR-10 as in Section 5.2.

Table 4. FID vs. NFE on AFHQv2.

NFE	3	9	13	19	35
EDM	375.76	27.88	7.56	2.99	2.11
ART-RL	<b>243.48</b>	<b>20.48</b>	<b>6.12</b>	<b>2.85</b>	<b>2.10</b>

**AFHQv2.** On AFHQv2, ART-RL attains lower FID than EDM at every NFE, with particularly large gains at small and moderate budgets and a consistent, albeit smaller, advantage even at the highest NFEs.

Table 5. FID vs. NFE on FFHQ.

NFE	3	9	13	19	35
EDM	466.76	57.13	15.87	5.26	2.73
ART-RL	<b>305.97</b>	<b>35.73</b>	<b>11.08</b>	<b>4.31</b>	<b>2.67</b>

**FFHQ.** The FFHQ results show the same pattern, with ART-RL improving on EDM across the entire NFE range. The advantage is pronounced at small NFEs and remains visible even at the largest budget ( $\text{NFE} = 35$ ).

Table 6. FID vs. NFE on ImageNet.

NFE	3	9	13	19	35
EDM	437.42	35.32	8.18	3.68	<b>2.57</b>
ART-RL	<b>147.21</b>	<b>29.49</b>	<b>7.01</b>	<b>3.62</b>	<b>2.57</b>

**ImageNet.** On ImageNet, we observe the same qualitative trend: ART-RL improves clearly over EDM at low and mid NFEs, with the gains narrowing at larger budgets.

**Qualitative results.** We provide qualitative comparisons for AFHQv2 and FFHQ in the appendix. For ImageNet,

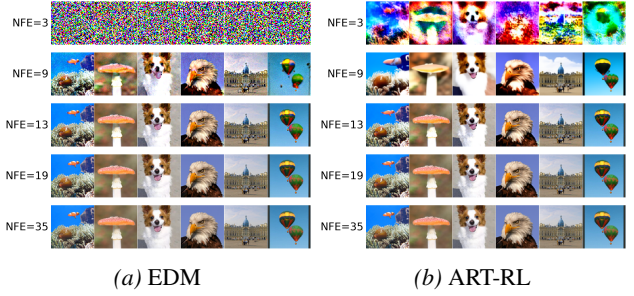


Figure 2. ImageNet samples under EDM and ART-RL schedules at increasing NFEs (top to bottom).

Figure 2 shows a representative example. At  $\text{NFE} = 3$ , ART-RL already produces partially recognizable objects, whereas the EDM outputs remain close to noise. At  $\text{NFE} = 9$ , ART-RL produces more coherent global structure with a smoother appearance, while EDM is still visibly noisy. For larger NFEs, both schedules produce high-quality images and the visual differences are subtle, consistent with the small FID gaps.

## 6. Conclusions

We introduce ART for timestep allocation in score-based diffusion models, along with ART-RL, a reinforcement learning algorithm that provides a continuous-time optimal control perspective instead of ad hoc schedule design. ART treats the time-warping rate as a learnable control and redistributes computation along the reverse trajectory under a fixed time budget. ART-RL uses an auxiliary, randomized formulation that leads to a numerically tractable actor-critic method, yielding a theoretically grounded schedule learned in a data-driven manner. Empirically, ART-RL improves sample quality given a same evaluation budget, generalizes across step counts and datasets and, once trained on CIFAR-10, its distilled time-only schedule serves as a training-free drop-in replacement in existing pipelines such as EDM on AFHQv2, FFHQ, and ImageNet.

While our approach has strong empirical performance and supporting theory, several directions await. Our analysis focuses on probability flow ODEs, yet extending ART to stochastic samplers may lead to different time-allocation behaviors. The learning objective relies on a surrogate for the Euler local truncation error; alternative surrogates or higher-order, solver-aware (e.g. Heun) criteria may better align the control formulation with practical integrators. Distillation collapses the learned policy to a time-only schedule, which works well empirically, but it remains unclear when state dependence matters. More broadly, adaptive time reparameterization for diffusion sampling is still nascent, and ART with ART-RL offers a control-theoretic starting point for systematic schedule design in generative diffusion models.



## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

- Benton, J., De Bortoli, V., Doucet, A., and Deligiannidis, G. Nearly  $d$ -linear convergence bounds for diffusion models via stochastic localization. In *ICLR*, 2024.
- Chen, H., Lee, H., and Lu, J. Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *ICML*, pp. 4735–4763, 2023a.
- Chen, S., Chewi, S., Li, J., Li, Y., Salim, A., and Zhang, A. R. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. In *ICLR*, 2023b.
- Choi, Y., Uh, Y., Yoo, J., and Ha, J.-W. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- Gao, X., Zha, J., and Zhou, X. Y. Reward-directed score-based diffusion models via q-learning. 2024. arXiv:2409.04832, To appear in *J. Mach. Learn. Res.*
- Google. State-of-the-art video and image generation with veo 2 and imagen 3. <https://blog.google/technology/google-labs/video-image-generation-update-december-2024/>, 2024. Accessed: 2025-09-17.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. In *Neurips*, volume 33, pp. 6840–6851, 2020.
- Huang, D. Z., Huang, J., and Lin, Z. Convergence analysis of probability flow ode for score-based generative models. 2025a. arXiv:2404.09730. To appear in *IEEE Trans. Inf. Theory*.
- Huang, Y. and Zhou, X. Y. Data-driven exploration for a class of continuous-time linear-quadratic reinforcement learning problems. 2025. arXiv:2507.00358.
- Huang, Y., Jia, Y., and Zhou, X. Achieving mean-variance efficiency by continuous-time reinforcement learning. In *Proceedings of the Third ACM International Conference on AI in Finance*, pp. 377–385, 2022.
- Huang, Y., Jia, Y., and Zhou, X. Y. Mean-variance portfolio selection by continuous-time reinforcement learning: Algorithms, regret analysis, and empirical study. 2024. arXiv:2412.16175.
- Huang, Y., Jia, Y., and Zhou, X. Y. Sublinear regret for a class of continuous-time linear-quadratic reinforcement learning problems. *SIAM Journal on Control and Optimization*, 63(5):3452–3474, 2025b.
- Jia, Y. and Zhou, X. Y. Policy evaluation and temporal-difference learning in continuous time and space: A martingale approach. *J. Mach. Learn. Res.*, 23(154):1–55, 2022a.
- Jia, Y. and Zhou, X. Y. Policy gradient and actor-critic learning in continuous time and space: Theory and algorithms. *J. Mach. Learn. Res.*, 23(154):1–55, 2022b.
- Jia, Y. and Zhou, X. Y.  $q$ -learning in continuous time. *J. Mach. Learn. Res.*, 24(161):1–61, 2023.
- Karras, T., Laine, S., and Aila, T. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4401–4410, 2019.
- Karras, T., Aittala, M., Aila, T., and Laine, S. Elucidating the design space of diffusion-based generative models. In *Neurips*, volume 35, pp. 26565–26577, 2022.
- Khanna, S., Kharbanda, S., Li, S., Varma, H., Wang, E., Birnbaum, S., Luo, Z., Miraoui, Y., Palrecha, A., and Ermon, S. Mercury: Ultra-fast language models based on diffusion. 2025. arXiv:2506.17298.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. 2009.
- Lee, H., Lu, J., and Tan, Y. Convergence for score-based generative modeling with polynomial complexity. In *Neurips*, volume 35, pp. 22870–22882, 2022.
- Li, G. and Yan, Y. Adapting to unknown low-dimensional structures in score-based diffusion models. In *Neurips*, volume 37, pp. 126297–126331, 2024.
- Li, G., Wei, Y., Chen, Y., and Chi, Y. Towards faster non-asymptotic convergence for diffusion-based generative models. In *ICLR*, 2024.
- Nie, S., Zhu, F., You, Z., Zhang, X., Ou, J., Hu, J., Zhou, J., Lin, Y., Wen, J.-R., and Li, C. Large language diffusion models. 2025. arXiv:2502.09992.
- OpenAI. Sora: Creating video from text. <https://openai.com/sora>, 2024. Accessed: 2025-09-17.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. Hierarchical text-conditional image generation with clip latents. 2022. arXiv:2204.06125.

- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *CVPR*, pp. 10684–10695, 2022.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., and Bernstein, M. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- Singer, U., Polyak, A., Hayes, T., Yin, X., An, J., Zhang, S., Hu, Q., Yang, H., Ashual, O., and Gafni, O. Make-a-video: Text-to-video generation without text-video data. In *ICLR*, 2023.
- Song, J., Meng, C., and Ermon, S. Denoising diffusion implicit models. In *ICLR*, 2021a.
- Song, Y. and Ermon, S. Generative modeling by estimating gradients of the data distribution. In *Neurips*, volume 32, pp. 11918–11930, 2019.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021b.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- Tang, W. and Zhao, H. Score-based diffusion models via stochastic differential equations. *Statistics Surveys*, 19: 28–64, 2025.
- Tang, W. and Zhou, X. Y. Regret of exploratory policy improvement and  $q$ -learning. 2024. arXiv:2411.01302.
- Wang, H., Zariphopoulou, T., and Zhou, X. Y. Reinforcement learning in continuous time and space: A stochastic control approach. *Journal of Machine Learning Research*, 21(198):1–34, 2020.
- Wu, Y., Chen, Y., and Wei, Y. Stochastic Runge-Kutta methods: Provable acceleration of diffusion models. 2024. arXiv:2410.04760.
- Zhang, Q. and Chen, Y. Fast sampling of diffusion models with exponential integrator. In *ICLR*, 2023.
- Zhang, Q., Song, J., and Chen, Y. Improved order analysis and design of exponential integrator for diffusion models sampling. 2023. arXiv:2308.02157.
- Zhao, H., Tang, W., and Yao, D. D. Policy optimization for continuous reinforcement learning. In *Neurips*, volume 36, 2023.
- Zhao, H., Chen, H., Zhang, J., Yao, D., and Tang, W. Scores as Actions: a framework of fine-tuning diffusion models by continuous-time reinforcement learning. 2024. arXiv:2409.08400.
- Zhao, H., Chen, H., Zhang, J., Yao, D., and Tang, W. Score as Action: Fine tuning diffusion generative models by continuous-time reinforcement learning. In *ICML*, 2025.
- Ziebart, B. D., Maas, A. L., Bagnell, J. A., and Dey, A. K. Maximum entropy inverse reinforcement learning. In *Aaai*, volume 8, pp. 1433–1438. Chicago, IL, USA, 2008.

## A. Proofs for Theorems 3.1 and 3.2

**Proof of Theorem 3.1.** We start from the definition  $V^{(\lambda)}(t, x, \psi) = V(t, x, \psi) + \lambda t$ . Differentiating gives

$$V_t^{(\lambda)}(t, x, \psi) = V_t(t, x, \psi) + \lambda, \quad V_x^{(\lambda)}(t, x, \psi) = V_x(t, x, \psi), \quad V_\psi^{(\lambda)}(t, x, \psi) = V_\psi(t, x, \psi).$$

Next, simplify the auxiliary HJB (16). Since

$$-|Q(x, \psi)| \left( \mu^2 + \frac{\lambda}{|Q(x, \psi)|} \right) = -|Q(x, \psi)| \mu^2 - \lambda,$$

equation (16) is equivalent to

$$V_t^{(\lambda)} + \sup_{\mu} \left\{ (V_x^{(\lambda)\top} F(x, \psi) + V_\psi^{(\lambda)} - \gamma) \mu - |Q(x, \psi)| \mu^2 - \lambda \right\} = 0. \quad (24)$$

Substitute the derivatives of  $V^{(\lambda)}$  into (24):

$$(V_t + \lambda) + \sup_{\mu} \left\{ (V_x^\top F(x, \psi) + V_\psi - \gamma) \mu - |Q(x, \psi)| \mu^2 - \lambda \right\} = 0.$$

The  $\lambda$  terms cancel, and we obtain

$$V_t + \sup_{\mu} \left\{ (V_x^\top F(x, \psi) + V_\psi - \gamma) \mu - |Q(x, \psi)| \mu^2 \right\} = 0.$$

This is exactly the original HJB (15) after renaming the maximization variable from  $\theta$  to  $\mu$ . Therefore, if  $V$  satisfies (15), then  $V^{(\lambda)}$  satisfies (16).

Finally, check the terminal condition. Using  $V(T, x, \psi) = \gamma T$ ,

$$V^{(\lambda)}(T, x, \psi) = V(T, x, \psi) + \lambda T = \gamma T + \lambda T = (\gamma + \lambda)T,$$

which matches the terminal condition of (16). Hence  $V^{(\lambda)}$  is a classical solution of (16).  $\square$

**Proof of Theorem 3.2.** We split the proof into three steps: (i) compute the maximizer of the HJB, (ii) prove a verification inequality  $V^{(\lambda)} \geq J^{\pi^{(\lambda)}}$  for any admissible policy, and (iii) show equality for the specific Gaussian policy  $\pi^{(\lambda)*}$ .

**Step 1: maximizer of the quadratic Hamiltonian.** Fix  $(t, x, \psi)$ . Consider the function of the scalar decision variable  $\mu$ :

$$\mathcal{H}(\mu) := (V_x^{(\lambda)\top} F(x, \psi) + V_\psi^{(\lambda)} - \gamma) \mu - |Q(x, \psi)| \mu^2.$$

This is a concave quadratic in  $\mu$  because  $|Q(x, \psi)| \geq 0$ . Differentiating with respect to  $\mu$  and setting to zero gives

$$\frac{d}{d\mu} \mathcal{H}(\mu) = V_x^{(\lambda)\top} F(x, \psi) + V_\psi^{(\lambda)} - \gamma - 2|Q(x, \psi)| \mu = 0,$$

hence the maximizer is

$$\mu^*(t, x, \psi) = \frac{V_x^{(\lambda)\top} F(x, \psi) + V_\psi^{(\lambda)} - \gamma}{2|Q(x, \psi)|}.$$

Using  $V_x^{(\lambda)} = V_x$  and  $V_\psi^{(\lambda)} = V_\psi$ , we obtain (18):

$$\mu^*(t, x, \psi) = \frac{V_x^\top F(x, \psi) + V_\psi - \gamma}{2|Q(x, \psi)|}.$$

In particular, the optimal mean control does not depend on  $\lambda$ .

**Step 2: Verification inequality**  $V^{(\lambda)} \geq J^{\pi^{(\lambda)}}$  **for any admissible policy**  $\pi^{(\lambda)} \in \Pi^{(\lambda)}$ . Fix  $\lambda > 0$  and fix any admissible Gaussian policy  $\pi^{(\lambda)} \in \Pi^{(\lambda)}$ . Let  $\mu^{\pi^{(\lambda)}}(t, x, \psi)$  denote the mean of  $\pi^{(\lambda)}(\cdot | t, x, \psi)$ , so that

$$\pi^{(\lambda)}(\cdot | t, x, \psi) = \mathcal{N}\left(\mu^{\pi^{(\lambda)}}(t, x, \psi), \frac{\lambda}{|Q(x, \psi)|}\right).$$

Let  $(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))$  be the state trajectory under policy  $\pi^{(\lambda)}$ . Under the exploratory formulation, the induced controlled dynamics can be written in feedback form as

$$\dot{x}^{\pi^{(\lambda)}}(t) = \mu^{\pi^{(\lambda)}}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)) F(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)), \quad \dot{\psi}^{\pi^{(\lambda)}}(t) = \mu^{\pi^{(\lambda)}}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)).$$

Apply the chain rule to the term  $V^{(\lambda)}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))$ :

$$\frac{d}{dt} V^{(\lambda)}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)) = V_t^{(\lambda)} + V_x^{(\lambda)\top} \dot{x}^{\pi^{(\lambda)}}(t) + V_\psi^{(\lambda)} \dot{\psi}^{\pi^{(\lambda)}}(t).$$

Substituting the dynamics yields, for every  $t \in [s, T]$ ,

$$\frac{d}{dt} V^{(\lambda)}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)) = V_t^{(\lambda)} + (V_x^{(\lambda)\top} F(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)) + V_\psi^{(\lambda)}) \mu^{\pi^{(\lambda)}}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)).$$

Integrating from  $s$  to  $T$  gives

$$\begin{aligned} & V^{(\lambda)}(T, x^{\pi^{(\lambda)}}(T), \psi^{\pi^{(\lambda)}}(T)) - V^{(\lambda)}(s, y, \phi) \\ &= \int_s^T \left[ V_t^{(\lambda)} + (V_x^{(\lambda)\top} F(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)) + V_\psi^{(\lambda)}) \mu^{\pi^{(\lambda)}}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)) \right] dt, \end{aligned} \quad (25)$$

where we condition on  $x^{\pi^{(\lambda)}}(s) = y$  and  $\psi^{\pi^{(\lambda)}}(s) = \phi$ .

Now add the running cost of the auxiliary objective to both sides of (25):

$$\int_s^T \left( -|Q(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))| [\mu^{\pi^{(\lambda)}}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))]^2 - \lambda - \gamma \mu^{\pi^{(\lambda)}}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)) \right) dt.$$

We obtain the identity

$$\begin{aligned} & V^{(\lambda)}(T, x^{\pi^{(\lambda)}}(T), \psi^{\pi^{(\lambda)}}(T)) - V^{(\lambda)}(s, y, \phi) + \int_s^T \left( -|Q| [\mu^{\pi^{(\lambda)}}]^2 - \lambda - \gamma \mu^{\pi^{(\lambda)}} \right) dt \\ &= \int_s^T \left[ V_t^{(\lambda)} + (V_x^{(\lambda)\top} F + V_\psi^{(\lambda)} - \gamma) \mu^{\pi^{(\lambda)}} - |Q| [\mu^{\pi^{(\lambda)}}]^2 - \lambda \right] dt, \end{aligned} \quad (26)$$

where, for readability, all terms  $|Q|$ ,  $F$ ,  $V_t^{(\lambda)}$ ,  $V_x^{(\lambda)}$ ,  $V_\psi^{(\lambda)}$ , and  $\mu^{\pi^{(\lambda)}}$  inside the integral are evaluated at  $(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))$ .

Define, for each  $(t, x, \psi)$ , the scalar function

$$\mathcal{G}(t, x, \psi; u) := (V_x^{(\lambda)\top} F(x, \psi) + V_\psi^{(\lambda)} - \gamma)u - |Q(x, \psi)| u^2 - \lambda.$$

Then the auxiliary HJB (16) is equivalent to

$$V_t^{(\lambda)}(t, x, \psi) + \sup_u \mathcal{G}(t, x, \psi; u) = 0.$$

Therefore, for any particular choice  $u = \mu^{\pi^{(\lambda)}}(t, x, \psi)$ ,

$$V_t^{(\lambda)}(t, x, \psi) + \mathcal{G}(t, x, \psi; \mu^{\pi^{(\lambda)}}(t, x, \psi)) \leq 0.$$



Applying this pointwise along the trajectory  $(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))$  implies that the integrand in (26) is nonpositive for all  $t \in [s, T]$ , hence

$$V^{(\lambda)}(T, x^{\pi^{(\lambda)}}(T), \psi^{\pi^{(\lambda)}}(T)) - V^{(\lambda)}(s, y, \phi) + \int_s^T \left( -|Q| [\mu^{\pi^{(\lambda)}}]^2 - \lambda - \gamma \mu^{\pi^{(\lambda)}} \right) dt \leq 0.$$

Rearranging and taking conditional expectation given  $x^{\pi^{(\lambda)}}(s) = y, \psi^{\pi^{(\lambda)}}(s) = \phi$  yields

$$\begin{aligned} V^{(\lambda)}(s, y, \phi) \geq \mathbb{E} \left[ \int_s^T \left( -|Q(x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))| [\mu^{\pi^{(\lambda)}}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t))]^2 - \lambda - \gamma \mu^{\pi^{(\lambda)}}(t, x^{\pi^{(\lambda)}}(t), \psi^{\pi^{(\lambda)}}(t)) \right) dt \right. \\ \left. + V^{(\lambda)}(T, x^{\pi^{(\lambda)}}(T), \psi^{\pi^{(\lambda)}}(T)) \mid x^{\pi^{(\lambda)}}(s) = y, \psi^{\pi^{(\lambda)}}(s) = \phi \right]. \end{aligned} \quad (27)$$

By definition of  $J^{\pi^{(\lambda)}}$ , the right-hand side is exactly  $J^{\pi^{(\lambda)}}(s, y, \phi)$ . This proves  $V^{(\lambda)}(s, y, \phi) \geq J^{\pi^{(\lambda)}}(s, y, \phi)$  for all admissible  $\pi^{(\lambda)} \in \Pi^{(\lambda)}$ .

**Step 3: Achieving equality and concluding optimality.** Equality in (27) holds if the policy mean  $\mu^{\pi^{(\lambda)}}(t, x, \psi)$  attains the supremum in the HJB, that is, if for every  $(t, x, \psi)$ ,

$$\mu^{\pi^{(\lambda)}}(t, x, \psi) \in \arg \max_u \mathcal{G}(t, x, \psi; u).$$

Since  $\mathcal{G}(t, x, \psi; u)$  is a concave quadratic in  $u$ , the maximizer is unique and given by

$$\mu^*(t, x, \psi) = \frac{V_x^{(\lambda)\top} F(x, \psi) + V_\psi^{(\lambda)} - \gamma}{2|Q(x, \psi)|}.$$

Using Theorem 3.1, we have  $V_x^{(\lambda)} = V_x$  and  $V_\psi^{(\lambda)} = V_\psi$ , hence

$$\mu^*(t, x, \psi) = \frac{V_x^\top F(x, \psi) + V_\psi - \gamma}{2|Q(x, \psi)|}.$$

Therefore, the Gaussian policy  $\pi^{(\lambda)*}$  defined by

$$\pi^{(\lambda)*}(\cdot \mid t, x, \psi) = \mathcal{N}\left(\mu^*(t, x, \psi), \frac{\lambda}{|Q(x, \psi)|}\right)$$

achieves equality in (27), which implies

$$V^{(\lambda)}(s, y, \phi) = J^{\pi^{(\lambda)*}}(s, y, \phi) = \sup_{\pi^{(\lambda)} \in \Pi^{(\lambda)}} J^{\pi^{(\lambda)}}(s, y, \phi).$$

Hence  $\pi^{(\lambda)*}$  is an optimal policy for the auxiliary problem.

Finally, since  $\mu^*(t, x, \psi)$  does not depend on  $\lambda$ , setting  $\lambda = 0$  yields the original ART control problem, and the same verification argument shows that  $\mu^*$  is an optimal policy for (10).  $\square$

## B. Additional Numerical Results

### B.1. Results for One-Dimensional Study

Figure 3 shows the empirical mean of the executed control  $\theta$  together with the 99 percent confidence band computed from the last 10,000 trajectories in the one-dimensional experiment. As in the main text, each trajectory is normalized so that the induced terminal time satisfies  $\psi(T) = T$ . The confidence band is extremely narrow and visually indistinguishable from the mean curve, confirming that in this setting the learned control exhibits negligible variability across trajectories and can be treated as an effectively deterministic function of time.

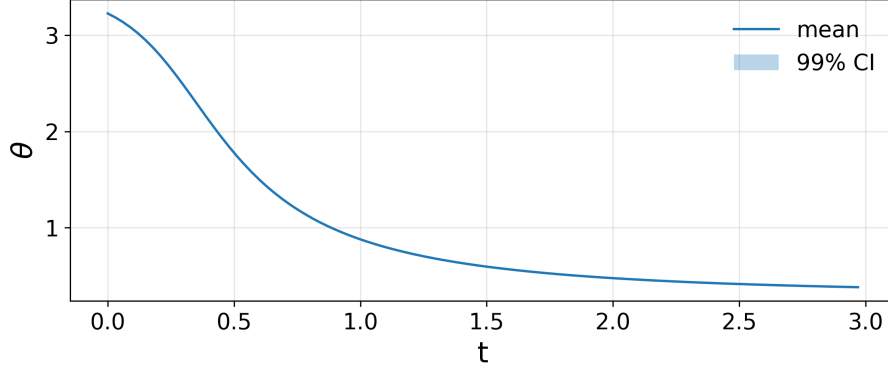


Figure 3. Empirical mean of the executed control  $\theta$  and its 99 percent confidence interval, based on the last 10,000 trajectories in the one-dimensional experiment.

Table 7. CIFAR-10 sampling with Euler updates in the EDM pipeline. We compare Uniform, EDM, and ART-RL time grids across representative step counts. Under Euler updates,  $\text{NFE} = K$ . Lower is better.

NFE	2	5	7	10	30	50	80
Uniform	280.504	214.599	194.398	174.146	85.8266	53.3991	34.9892
EDM	295.65	49.0991	27.7335	15.5683	4.21215	3.00838	2.50025
ART-RL	<b>109.11</b>	<b>28.1562</b>	<b>23.8837</b>	<b>14.341</b>	<b>4.0591</b>	<b>2.94321</b>	<b>2.4605</b>

## B.2. CIFAR-10 Euler Ablation under the EDM Pipeline

Our training objective is theoretically motivated by an Euler local discretization error proxy; so using Euler updates is a natural way to validate whether the learned time grid improves sampling under the same proxy. This ablation is not meant to suggest that Euler is the preferred solver for image generation in general; rather, it serves as a controlled check that aligns the evaluation solver with the proxy used in training.

We rerun CIFAR-10 using Euler updates within the official EDM pipeline, while keeping all other components unchanged. Under Euler updates, each step requires one score evaluation, so the step count equals the number of function evaluations, that is,  $\text{NFE} = K$ . We compare three time grids, Uniform, EDM, and ART-RL, across  $K \in \{2, 5, 7, 10, 30, 50, 80\}$ . As shown by Table 7, at all tested  $K$ , ART-RL consistently outperforms both Uniform and EDM grids.

## B.3. Qualitative Results for CIFAR-10

For completeness, we include additional visual results for CIFAR-10 at all NFEs considered in the main text, as shown in Figure 4. Each panel displays a grid of samples generated under a timestep schedule (Uniform, EDM, or ART-RL). Within each panel, the rows correspond to increasing NFEs, allowing visual inspection of how sample quality improves as more function evaluations are used. The ART-RL samples exhibit faster refinement across NFEs, consistent with the quantitative FID results reported in Section 5.2.

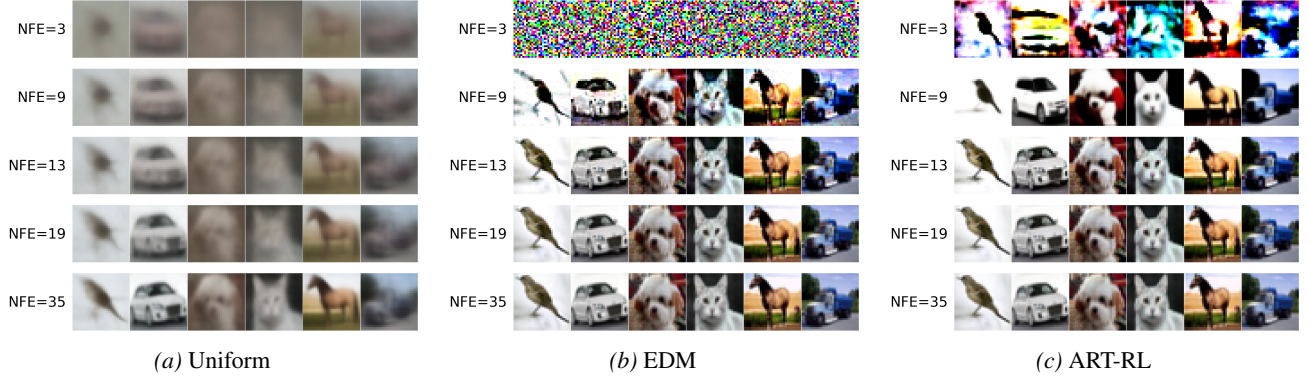


Figure 4. CIFAR-10 samples across timesteps for the three schedules (Uniform, EDM, ART-RL). Each panel shows a grid where rows correspond to increasing NFEs.

#### B.4. Qualitative Results for the Generalization of the ART-RL Time Schedules

This appendix provides visual samples for the experiments in Section 5.3. For the CIFAR-10 interpolation and extrapolation study (Section 5.3.1), and for the cross-dataset transfer experiments on AFHQv2, FFHQ, and ImageNet (Section 5.3.2), we display grids of generated images for EDM and ART-RL. Each panel shows samples at increasing NFEs across rows, complementing the quantitative comparisons in the main text.

##### B.4.1. CIFAR-10: INTERPOLATED AND EXTRAPOLATED TIME GRIDS

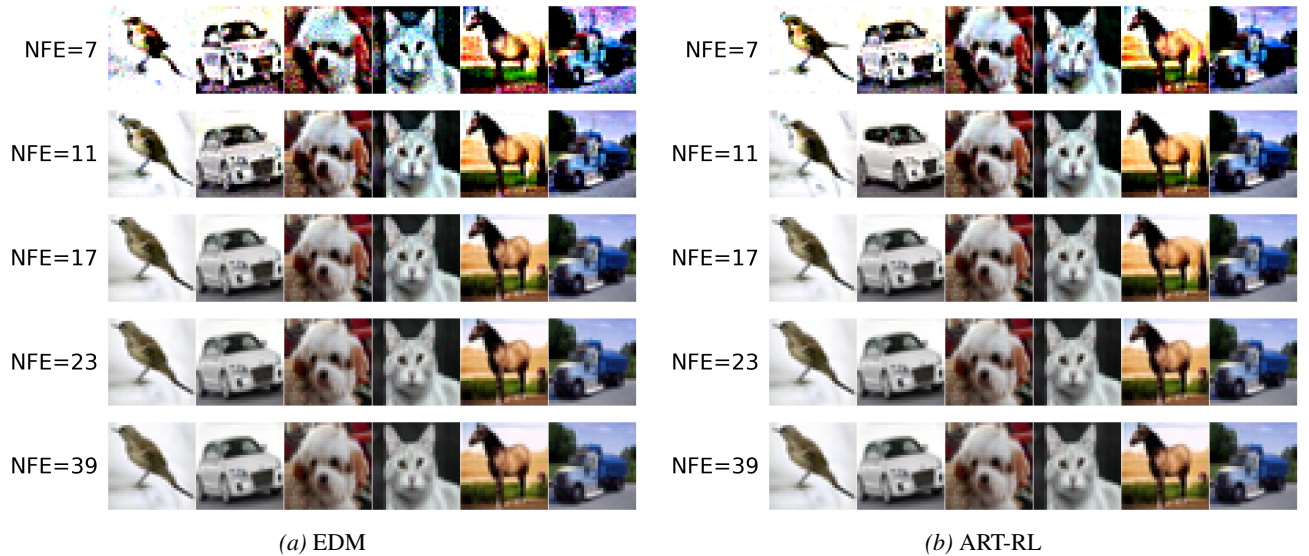


Figure 5. CIFAR-10 samples across timesteps for interpolated and extrapolated grids (EDM and ART-RL). Each panel shows a grid where rows correspond to increasing NFEs.

## B.4.2. AFHQv2

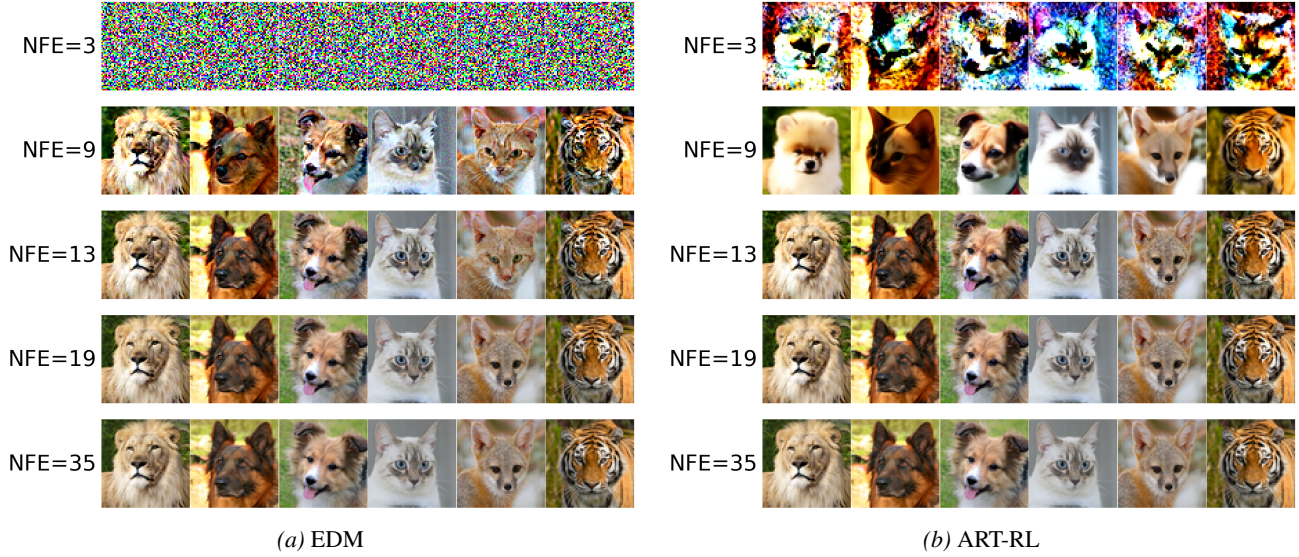


Figure 6. AFHQv2 samples across timesteps for the two schedules (EDM and ART-RL). Each panel shows a grid where rows correspond to increasing NFEs.

## B.4.3. FFHQ

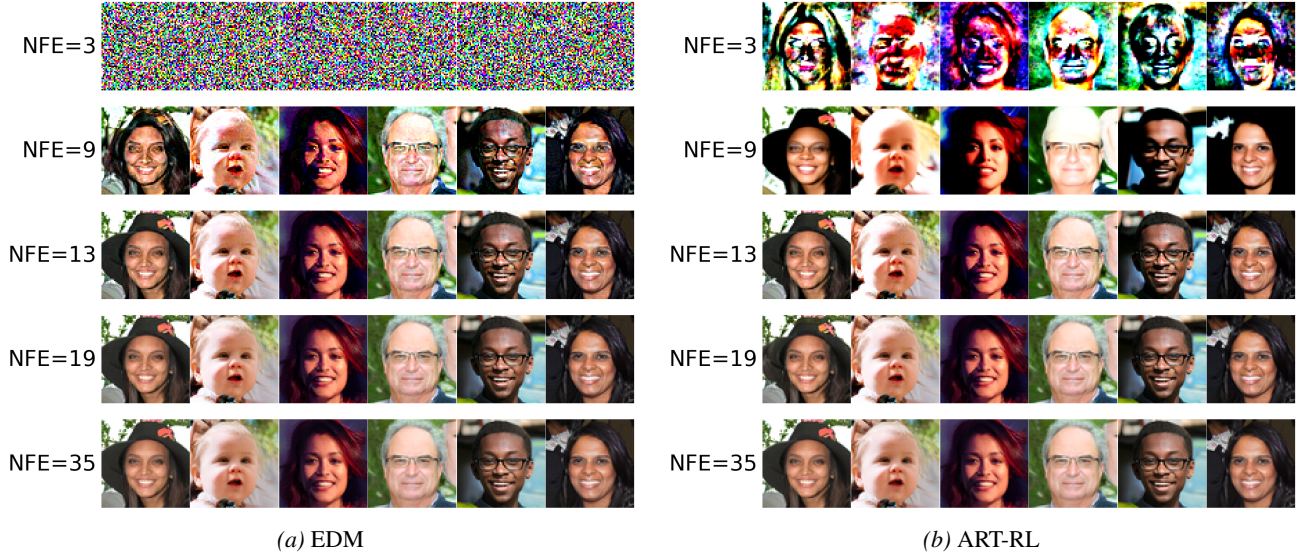


Figure 7. FFHQ samples across timesteps for the two schedules (EDM and ART-RL). Each panel shows a grid where rows correspond to increasing NFEs.

## C. Reproducibility and Training Overhead

- Our image experiments follow the official EDM pipeline and keep the score model, solver, noise-conditioning, and EDM hyperparameters fixed. ART-RL replaces only the time grid. The EDM schedule uses the standard exponent  $\rho = 7$  in all image experiments. In the one-dimensional example, we additionally report a tuned variant with  $\rho = 3$ .
- In Algorithm 1, we discretize the new clock  $t$  on a uniform grid  $0 = t_0 < t_1 < \dots < t_K = T$  with  $\Delta t = T/K$ . Each RL iteration rolls out one backward trajectory under the current policy and performs one critic update and one actor



update using the Riemann discretized moment conditions.

- The actor and critic are 3-layer MLPs with hidden width 128 and Softplus activations. For images, the networks do not process the raw tensor  $x$  directly. Instead,  $x$  is represented by a low-dimensional feature vector computed from quantities already evaluated along the rollout, including  $t, \psi$ .
- The Gaussian policy variance uses the parameterization in the paper,

$$\hat{\pi}_{\vartheta_a}(\cdot \mid t, x, \psi) = \mathcal{N}\left(\mu(t, x, \psi), \frac{\lambda}{|Q(x, \psi)| \vee \varepsilon}\right),$$

with  $\lambda = 10^{-1}$  and  $\varepsilon = 10^{-6}$ .

- We run  $N = 5,000$  iterations. We use Adam for both actor and critic with learning rate  $10^{-4}$ ,  $(\beta_1, \beta_2) = (0.9, 0.999)$ , and no weight decay. The Lagrange multiplier is updated by stochastic approximation with step size  $10^{-4}$ .
- Computing  $Q(x, \psi)$  (Eq. (8)) requires  $\nabla_x \hat{S}(s, x)$  and  $\partial_s \hat{S}(s, x)$ , where  $s = T - \psi$ . In implementation,  $\nabla_x \hat{S}(s, x)$  is never formed explicitly. We obtain the required quantities through automatic differentiation using Jacobian vector products, together with differentiation of the score output with respect to the scalar time input  $s$ . Per training step, this adds two derivative queries, one Jacobian vector product and one time derivative query, on top of one score evaluation.
- To ensure the score model is always queried at a valid noise level, we enforce  $\psi(t) \in [0, T]$ , hence  $s = T - \psi(t) \in [0, T]$ , along executed trajectories. In particular, we normalize the realized distilled  $\theta$  sequence on each trajectory so that the induced total time change matches the constraint  $\sum_{k=0}^{K-1} \theta(t_k) \Delta t = T$ , which removes numerical over or under shoot of  $\psi$ .
- After distillation, sampling uses only a fixed precomputed time grid. The actor, critic, and  $Q(x, \psi)$  are not evaluated at inference. As a result, the per step sampling runtime matches EDM and Uniform under the same solver and score model.