

---

# FASTER : Value-Guided Sampling for Fast RL

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Some of the most performant reinforcement learning algorithms today can be  
2 prohibitively expensive as they use test-time scaling methods such as sampling  
3 multiple action candidates and selecting the best one. In this work, we propose  
4 **FASTER**, a method for getting the benefits of sampling-based test-time scaling of  
5 diffusion-based policies without the computational cost by tracing the performance  
6 gain of action samples back to earlier in the denoising process. Our key insight is  
7 that we can model the denoising of multiple action candidates and selecting the best  
8 one as a Markov Decision Process (MDP) where the goal is to progressively filter  
9 action candidates before denoising is complete. With this MDP, we can learn a  
10 policy and value function in the denoising space that predicts the downstream value  
11 of action candidates in the denoising process and filters them while maximizing  
12 returns. The result is a method that is lightweight and can be plugged into existing  
13 generative RL algorithms. Across challenging long-horizon manipulation tasks in  
14 online and batch-online RL, **FASTER** consistently improves the underlying policies  
15 and achieves the best overall performance among the compared methods. Applied  
16 to a pretrained VLA, **FASTER** achieves the same performance while substantially  
17 reducing training and inference compute requirements.

18 Code: <https://anonymous.4open.science/r/faster-C250>

## 19 1 Introduction

20 Recent reinforcement learning (RL) methods have demonstrated strong performance using expressive  
21 policy backbones such as diffusion models which are widely used in domains such as image/video  
22 generation and robotics. However, some of the most performant of these algorithms are prohibitively  
23 expensive both during training and at test time. This is because many of these algorithms rely on  
24 sampling multiple candidate actions and selecting the best one, for example as determined by which  
25 one has the highest value. While effective, such approaches incur substantial computational costs.  
26 This is especially problematic for large models like modern Vision-Language-Action (VLA) models,  
27 which already operate near the limits of acceptable latency. A modest sample count of eight will  
28 multiply inference costs accordingly and can render these methods impractical in latency-sensitive or  
29 resource-constrained settings. In this work, we ask: can we recover the benefits of sampling-based  
30 test-time scaling without suffering its computational cost?

31 This high computational cost comes from policies with generative models such as diffusion or flow  
32 matching that need to denoise all candidates. Prior works have shown that test-time scaling methods  
33 for generative models are performant. Methods such as best-of- $N$  sampling and self-consistency  
34 sample multiple candidates and select the one that maximizes a given metric. While these methods  
35 are effective at maximizing performance, they scale poorly and require full denoising over all  
36 action candidates. Distillation-based approaches amortize this cost by training the policy to directly  
37 reproduce high-value behaviors, but require training a separate policy which can be expensive. The  
38 core difficulty underlying all of these approaches is that, without observing the fully denoised samples,  
39 it is hard to know which ones are worth executing since there is no direct supervision linking a given  
40 noise sample to its downstream outcome, and the mapping from noise to action is only defined  
41 implicitly through the denoising process.

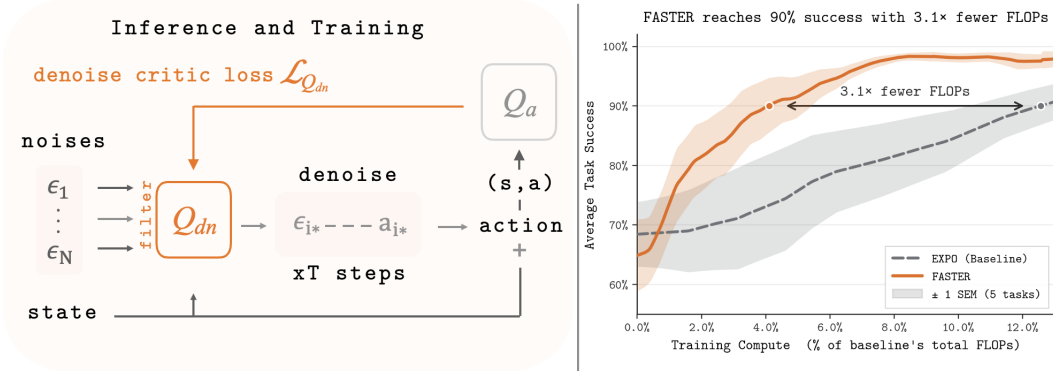


Figure 1: **Left: Overview of FASTER.** Instead of denoising all  $N$  candidates and selecting the best action post-hoc (best-of- $N$ ), FASTER learns a denoise critic  $Q^{dn}$  that scores action samples *during* denoising, often directly on the initial noise itself. **Right: Performance of FASTER compared with baseline** averaged across tasks for  $\pi_{0.5}$ . The FLOPs are normalized on the x-axis against the amount of FLOPs needed for the baseline to converge.

42 In this work, our key insight is that we can model the denoising of multiple candidate samples and  
 43 the selection of the best one as a Markov Decision Process (MDP) where the goal is to progressively  
 44 filter candidates before denoising is complete. With this MDP, we can learn a denoise Q-function and  
 45 policy with traditional temporal difference learning that decides which actions to keep and remove  
 46 while maximizing the returns. To supervise the value function without requiring exhaustive rollouts,  
 47 we use the value function over final actions as the reward of the action filtering MDP, enabling the  
 48 model to propagate outcome information back to the denoising space without direct outcome labels  
 49 for every candidate. Practically, we find that filtering candidates at the noise level—which is the  
 50 cheapest instantiation computationally—achieves performance equivalent to fully denoising all action  
 51 samples and selecting the highest value action. The computational cost of denoising—which is the  
 52 dominant bottleneck at inference time—is thus reduced to that of a single rollout, regardless of the  
 53 number of action candidates. The result is a lightweight, modular approach that decouples the benefits  
 54 of broad action sample selection from the cost of full policy execution.

55 Our main contribution is a framework for obtaining the performance gains of sampling and selecting  
 56 the best of multiple actions of RL methods without the computational overhead of denoising all  
 57 actions. Unlike approaches that operate on actions, our method operates directly on the denoising  
 58 process to filter down candidates, making selection both cheap and modular. We evaluate our method  
 59 on challenging manipulation benchmarks, including Robomimic and LIBERO, across algorithms and  
 60 model capacities, demonstrating consistent gains over single-sample baselines and near-parity with  
 61 best-of- $N$  methods at a fraction of the computational overhead.

## 62 2 Related Work

63 **Test-time scaling and best-of- $N$  inference.** Significant work has shown test-time scaling methods  
 64 for generative models to be highly performant. These methods take several forms, including iterative  
 65 refinement, but the most common involve sampling multiple candidates and performing filtering or  
 66 refinement based on a group of samples. Best-of- $N$  sampling is the most standard of such methods,  
 67 with extensive use in many domains including in autoregressive language models [1, 2, 3], diffusion  
 68 models [4], and policy learning [5, 6, 7, 8].

69 **Diffusion policies and value-guided sampling.** Diffusion models have emerged as a strong policy  
 70 class for a wide variety of domains because they can represent rich, multimodal data distributions [5,  
 71 6, 9, 10, 11, 12, 13, 14]. Several recent works study how to improve these policies with value  
 72 functions or reduce their inference cost. Value-Guided Policy Steering (V-GPS) re-ranks actions  
 73 from frozen robot policies with a learned value function [15], RoboMonkey scales test-time sampling  
 74 and verification for vision-language-action policies [16], and EXPO couples an expressive base  
 75 policy with a lightweight edit policy and selects value-maximizing actions online [5]. Other work  
 76 amortizes or changes the policy itself: One-Step Diffusion Policy distills multi-step denoising into a  
 77 faster actor [17], while DSRL post-trains a behavior-cloned diffusion policy by running RL in its  
 78 latent-noise space [18]. Compared to these approaches, our method learns a separate critic over *which*

79 *actions to filter during denoising*, enabling candidate selection before actions are denoised. As such,  
80 our method can be directly applied to methods that can benefit from sampling multiple actions.

81 **Initial-noise and noise-space methods for diffusion models.** A rapidly growing body of literature  
82 shows that the initial seed of a diffusion model has a significant effect on generation quality. Prior  
83 work has shown that seed quality is linked to inversion stability [19, 20] and several methods explicitly  
84 optimize or learn better starting noises, including InitNO and FIND [21, 22], as well as plug-and-  
85 play noise refinement methods such as Golden Noise [23], NoiseRefine [24], and Noise-Level  
86 Guidance [25]. Noise Hypernetworks likewise seek to amortize test-time noise optimization into  
87 a learned auxiliary module [26]. More closely related, TTSnap [27] prunes candidate trajectories  
88 early by decoding intermediate estimates and scoring them with noise-aware reward models trained  
89 via self-distillation, allowing it to search over larger candidate sets. This is similar in spirit to our  
90 goal of recovering the benefits of best-of- $N$  without denoising every candidate, but the mechanism is  
91 distinct. TTSnap still performs partial denoising and reward evaluation at selected timesteps, and  
92 is studied in text-to-image generation with image-reward objectives from pretrained models. By  
93 contrast, we frame the problem of filtering action candidates as an MDP and learn a *state-conditioned*  
94 critic that predicts downstream value directly from the initial noise of a diffusion policy in the context  
95 of reinforcement learning.

### 96 3 Preliminaries

97 We consider problems modeled as Markov decision processes (MDPs), each defined by a tuple  
98  $\{\mathcal{S}, \mathcal{A}, r, \gamma, T, \rho\}$  where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward  
99 function,  $T(s'|s, a)$  comprises the transition dynamics,  $\gamma \in [0, 1]$  is the discount factor, and  $\rho(s)$   
100 is the initial state distribution. In this MDP, the RL agent observes state  $s_t$  at every timestep  $t$  and  
101 chooses action  $a_t$  from its policy  $\pi(a_t|s_t)$ . The tuple of states, actions, rewards, and next states  
102  $(s, a, r, s')$  is appended to a replay buffer  $D$  for training. The goal of RL is to maximize the expected  
103 sum of discounted returns  $\mathbb{E}_\pi[\sum_{t=0}^T \gamma^t r(s_t, a_t)]$ .

104 In this paper, we study the setting where the policy  $\pi$  is a diffusion or flow policy that requires  
105 denoising. We further focus on algorithms that perform test-time scaling through sampling multiple  
106 action candidates. Given  $N$  initial noise seeds  $\{\epsilon_i\}_{i=1}^N$ , the policy denoises them into  $N$  actions  
107  $\{a_i\}_{i=1}^N$ . The policy then selects the best action among the candidate set, for example by choosing  
108  $i^* = \arg \max_i Q(s, a_i)$  and executing  $a_{i^*}$ . Our main goal is to capture the performance gain from  
109 these sampling-based test-time scaling approaches without the computational cost of denoising all  $N$   
110 action candidates.

### 111 4 FASTER

112 To capture the benefits of test-time scaling without its computational overhead, we leverage a key  
113 insight: the sample variance that drives these gains can be traced to early in the denoising process.  
114 We introduce **FASTER**, a framework that models denoising as an MDP and learns a filtering policy to  
115 select promising actions before they are denoised. An overview is shown in [Figure 1](#). We first define  
116 the filtering MDP ([Section 4.1](#)), then describe policy learning within this MDP ([Section 4.2](#)). In  
117 practice, we find that sample variance is largely determined by the initial noise, motivating a practical  
118 instantiation that filters at the noise level — the earliest and computationally cheapest point in the  
119 denoising chain. We describe this and further practical implementations in [Section 4.3](#).

#### 120 4.1 Modeling Action Denoising as an MDP

121 We want to select the best action candidates *before* denoising completes. To do so, we model the  
122 reverse diffusion process as a filtering MDP over  $N$  noise candidates  $\{\epsilon_i\}_{i=1}^N$ ,  $\epsilon_i \sim \mathcal{N}(0, I)$ , where a  
123 policy progressively discards unpromising candidates at each denoising step. We illustrate this in  
124 [Figure 2](#).

125 **States.** A state  $\mathbf{s}_t = (s, t, \mathcal{C}_t, \{a_i^{(t)}\}_{i \in \mathcal{C}_t})$  consists of the environment state  $s$ , the denoising timestep  
126  $t \in \{T, \dots, 1\}$ , the surviving candidate set  $\mathcal{C}_t \subseteq \{1, \dots, N\}$ , and the partially denoised intermediates,  
127 from pure noise  $a_i^{(T)} = \epsilon_i$  to a clean action at  $t=1$ .

**Algorithm 1** Learning procedure in **FASTER**

**Require:**  $s, \pi_\theta, Q^{dn}, Q^a, N$   
 1: Sample  $\{\epsilon_i\}_{i=1}^N$  with  $\epsilon_i \sim \mathcal{N}(0, I)$   
 2:  $i^* \leftarrow \arg \max_i Q^{dn}(s, \epsilon_i)$   
 3:  $a_{i^*} \leftarrow \pi_\theta(s, \epsilon_{i^*})$   
 4:  $\mathcal{L}_{\text{reg}} = \|Q^{dn}(s, \epsilon_{i^*}) - \text{sg}[Q^a(s, a_{i^*})]\|_2^2$

**Algorithm 2** Inference procedure in **FASTER**

**Require:**  $s, \pi_\theta, Q^{dn}, N$   
 1: Sample  $\{\epsilon_i\}_{i=1}^N$  with  $\epsilon_i \sim \mathcal{N}(0, I)$   
 2:  $i^* \leftarrow \arg \max_i Q^{dn}(s, \epsilon_i)$   
 3:  $a_{i^*} \leftarrow \pi_\theta(s, \epsilon_{i^*})$   
 4: Execute  $a_{i^*}$

128 **Actions.** At each step the policy retains or discards each  
 129 candidate:  $m_{t,i} \in \{0, 1\}$  for  $i \in \mathcal{C}_t$ , with at least one retained  
 130 ( $\sum_i m_{t,i} \geq 1$ ). The surviving set becomes  $\mathcal{C}_{t-1} =$   
 131  $\{i : m_{t,i} = 1\}$ .

132 **Transitions.** A denoising step is applied to each candidate,  
 133 producing  $a_i^{(t-1)}$ . The episode ends when a single candi-  
 134 date remains ( $|\mathcal{C}_{t-1}| = 1$ ) or denoising finishes ( $t = 1$ ), at  
 135 which point the highest-scoring survivor is executed. The  
 136 episode may terminate at  $t > 1$ , in which case the remaining  
 137 candidate is denoised until  $t = 1$ .

138 **Reward.** Non-terminal steps receive zero reward. At  
 139 termination the reward is  $Q^a(s, a_{i^*})$ , where  $a_{i^*}$  is the  
 140 executed action.

141 **4.2 Learning the filtering policy**

142 To learn the filtering policy in the MDP above, we learn an  
 143 action-value function over filtering decisions. Concretely,  
 144 we learn the denoise critic  $Q^{dn}(s_t, \mathbf{m}_t)$ , where  $s_t = (s, t, \mathcal{C}_t, \{a_i^{(t)}\}_{i \in \mathcal{C}_t})$  is the filtering state and  
 145  $\mathbf{m}_t = \{m_{t,i}\}_{i \in \mathcal{C}_t}$  is the filtering action with  $m_{t,i} \in \{0, 1\}$  denoting whether candidate  $i$  is retained.  
 146 We train  $Q^{dn}$  with a temporal-difference (TD) objective.

$$\mathcal{L}_{\text{TD}} = \mathbb{E}_{(s_t, \mathbf{m}_t, s_{t-1})} \left[ \left( Q^{dn}(s_t, \mathbf{m}_t) - \left[ r_t + \gamma \max_{\mathbf{m}_{t-1}} \bar{Q}^{dn}(s_{t-1}, \mathbf{m}_{t-1}) \right] \right)^2 \right] \quad (1)$$

147 where  $r_t = 0$  for non-terminal transitions and, at termination,  $r_t = Q^a(s, a_{i^*})$  for the executed action  
 148  $a_{i^*}$  induced by the final surviving candidate.

149 At test time, we sample  $N$  noise candidates  $\{\epsilon_i\}_{i=1}^N$  with  $\epsilon_i \sim \mathcal{N}(0, I)$  and repeatedly apply the  
 150 filtering policy induced by  $Q^{dn}$  to choose filtering actions  $\mathbf{m}_t$  until a single candidate remains. Let  
 151  $i^*$  denote the surviving candidate index. That survivor is then denoised to obtain the final action  
 152  $a_{i^*} = \pi_\theta(s, \epsilon_{i^*})$ . Depending on the method (Section 4.3), we either use  $a_{i^*}$  as the final action or  
 153 consider this alongside a proposal action from an edit policy.

154 **4.3 Practical implementation**

155 In our setting, we restrict the MDP to only filter at  $t = T$ , reducing the multi-step filtering process  
 156 to a single decision. Since the filtering happens before any denoising, the state collapses to just the  
 157 environment state  $s$  (the denoising timestep is fixed at  $T$  and there is no history of prior filtering steps),  
 158 and the filtering ‘‘action’’ reduces to selecting a single candidate rather than producing a binary mask  
 159  $\mathbf{m}_t$  over survivors. Under this simplification, the joint  $Q^{dn}(s_t, \mathbf{m}_t)$  from Section 4.2 decomposes  
 160 into per-candidate scores: we learn a *noise-level critic*  $Q^{dn}(s, \epsilon)$  that maps an environment state  $s$   
 161 and a single noise sample  $\epsilon \sim \mathcal{N}(0, I)$  to the expected return of denoising that sample. We denote  
 162 the standard action-space critic as the *action-level critic*  $Q^a(s, a)$ , which provides the supervision  
 163 signal. We use a greedy filtering policy which selects  $i^* = \arg \max_i Q^{dn}(s, \epsilon_i)$  and denoises  $\epsilon_{i^*}$  into  
 164  $a_{i^*} = \pi_\theta(s, \epsilon_{i^*})$ .

165 Because the MDP is now one step, the TD objective in Section 4.2 simplifies to a one-step regression  
 166 against the action-level critic. We sample  $N$  noise candidates  $\{\epsilon_i\}_{i=1}^N$  with  $\epsilon_i \sim \mathcal{N}(0, I)$ , select

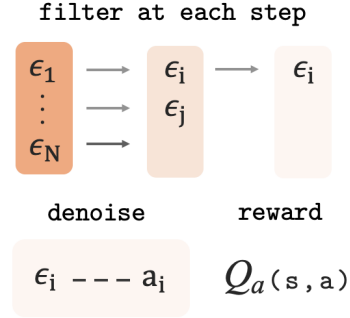
**Action Filtering MDP**

Figure 2: **Action Filtering MDP.** We model the process of denoising action candidates and selecting the best one as an MDP where the goal is to filter action samples during denoising while maximizing returns.

167 the highest-scoring index  $i^* = \arg \max_i Q^{dn}(s, \epsilon_i)$ , denoise  $\epsilon_{i^*}$  to obtain  $a_{i^*} = \pi_\theta(s, \epsilon_{i^*})$ , and  
 168 supervise  $Q^{dn}$  by regressing to the action-level critic:

$$\mathcal{L}_{\text{reg}} = \|Q^{dn}(s, \epsilon_{i^*}) - \text{sg}[Q^a(s, a_{i^*})]\|_2^2.$$

169 The regression-based learning and inference procedures are summarized in [Algorithms 1](#) and [2](#). At  
 170 test time, we select  $i^* = \arg \max_{i \in \{1, \dots, N\}} Q^{dn}(s, \epsilon_i)$  and denoise only the selected candidate  $\epsilon_{i^*}$  to  
 171 obtain the final action  $a_{i^*} = \pi_\theta(s, \epsilon_{i^*})$ .

172 We evaluate **FASTER** with two sampling-based methods in our experiments.

173 **FASTER-EXPO**. We implement **FASTER** on top of EXPO, a high-performing online RL method that  
 174 leverages a diffusion base policy combined with a lightweight edit policy. In this case, the selected  
 175 denoised action  $a_{i^*}$  serves as the base-policy proposal, and the downstream edit step follows EXPO  
 176 unchanged.

177 **FASTER-IDQL**. We implement **FASTER** on top of IDQL, which performs implicit action selection  
 178 from a diffusion policy without an edit policy. In this case, the denoised action  $a_{i^*}$  is executed  
 179 directly.

## 180 4.4 Computational profile

181 The key computational advantage of **FASTER** over best-of- $N$  sampling is that it replaces  $N$  full  
 182 denoising rollouts with  $N$  cheap noise-level critic evaluations followed by a single denoising rollout.  
 183 Let  $T$  denote the number of denoising steps,  $N$  the number of candidates, and  $F_{\text{actor}}, F_{\text{critic}}$  the  
 184 FLOPs for a single forward pass of the actor and critic, respectively. Standard best-of- $N$  sampling  
 185 denoises all  $N$  candidates for  $T$  steps each and then scores the resulting actions:

$$C_{\text{BoN}} = T \cdot N \cdot F_{\text{actor}} + N \cdot F_{\text{critic}}.$$

186 **FASTER** instead scores all  $N$  noise candidates with the noise critic and denoises only the top-scoring  
 187 one:

$$C_{\text{FASTER}} = N \cdot F_{\text{critic}} + T \cdot F_{\text{actor}}.$$

188 The savings come from reducing the actor cost from  $\mathcal{O}(TN)$  to  $\mathcal{O}(T)$ , eliminating  $(N-1)$  full  
 189 denoising passes. When  $T = N$  and  $F_{\text{actor}} = F_{\text{critic}} = F$ , the cost reduces from  $T(T+1) \cdot F$  to  
 190  $2T \cdot F$ —an approximate  $T/2 \times$  speedup.

191 In practice,  $T$  typically ranges from 5 to 20 denoising steps. For example, EXPO [\[5\]](#) uses  $T=10$   
 192 and  $N=8$ . While one-step diffusion methods [\[28, 29\]](#) attempt to distill multi-step generations into  
 193 a single-step process, state-of-the-art diffusion models in robotics, image generation, and video  
 194 generation [\[30, 31, 32, 33\]](#) continue to use multi-step sampling in this range, making the cost  
 195 reduction from **FASTER** broadly applicable.

196 In our primary experiments, the actor and critic share the same architecture and FLOP count ( $F_{\text{actor}} =$   
 197  $F_{\text{critic}}$ ), consistent with baselines such as [\[5\]](#) and recent work on large-scale RL [\[34\]](#). For our VLA  
 198 experiments, the asymmetry is even more favorable: the critic has 20M parameters versus the actor’s  
 199 3.3B, so  $F_{\text{critic}} \ll F_{\text{actor}}$  and the  $N$  noise-level evaluations add negligible overhead.

## 200 5 Experiments

201 The goal of our experiments is to answer the following core questions:

- 202 (Q1) How does **FASTER** perform compared to state-of-the-art RL methods in both the online setting  
 203 and the batch-online setting?
- 204 (Q2) Is **FASTER** able to recover the performance gains of test-time scaling while reducing the  
 205 computational cost?
- 206 (Q3) Can **FASTER** be applied to a pretrained Vision-Language-Action (VLA) model?
- 207 (Q4) What components of **FASTER** are most important for performance?

208 **Environments**. We evaluate **FASTER** on a set of 9 challenging manipulation tasks from Robomimic  
 209 and LIBERO. All tasks feature a sparse reward indicating task completion. These environments  
 210 involve controlling a 7-DoF robot arm to complete complex manipulation tasks. For Robomimic,  
 211 we evaluate on Lift, Can, Square, and Tool Hang, which require picking up a block, picking up a

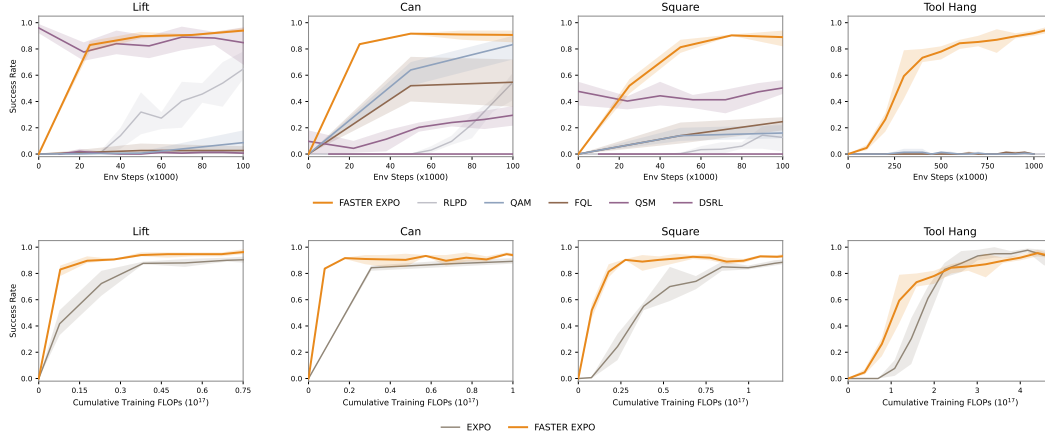


Figure 3: **Top:** Success rates of **FASTER** and baselines in the online settings. **FASTER-EXPO** outperforms strong baselines in sample efficiency. **Bottom:** Compute comparisons of **FASTER-EXPO** and EXPO. **FASTER** eliminates extra denoising during training and inference, yielding large FLOP reductions relative to the EXPO with comparable task performance.

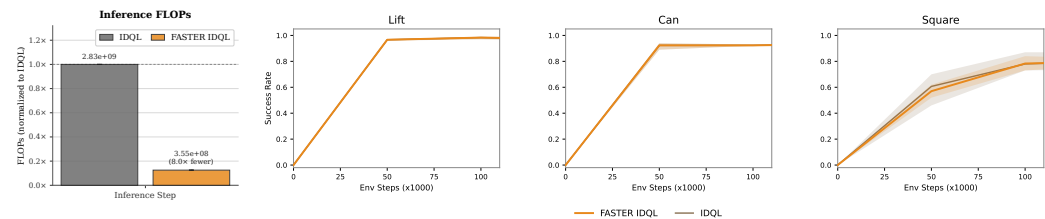


Figure 4: **Success rate and compute comparisons** of **FASTER-IDQL** and IDQL in the online setting. **FASTER** can be applied to IDQL to eliminate extra denoising rollouts at inference while obtaining the same performance in success rates (see Figure 3).

212 can and placing it, inserting a tool onto a square peg, and hanging a tool on a rack, respectively. For  
 213 LIBERO, we start from the pretrained pi05\_libero checkpoint from OpenPI [30], which is trained  
 214 on libero\_goal, libero\_object, libero\_spatial, and libero\_10 but not libero\_90; we  
 215 then select all held-out libero\_90 tasks on which the original policy achieves 40–60% success and  
 216 run RL on those tasks without any offline data.

217 **Baselines.** We compare our method to prior state-of-the-art methods for online and batch-online  
 218 reinforcement learning [35].

219 **EXPO [5].** EXPO is an online RL method that jointly learns an expressive diffusion policy alongside  
 220 a lightweight Gaussian edit policy that edits the actions sampled from the base policy toward a higher  
 221 value distribution.

222 **IDQL [6].** IDQL trains an expressive diffusion policy via imitation learning and uses implicit policy  
 223 extraction by performing best-of- $N$  sampling, selecting the action that maximizes the  $Q$ -value.

224 **RLPD [36].** RLPD is a highly sample-efficient algorithm that leverages prior data and samples  
 225 from it for learning. RLPD uses a simpler Gaussian policy and has been shown to outperform many  
 226 offline-to-online methods even without pretraining. For both evaluation settings, we run RLPD  
 227 without offline pretraining.

228 **QSM [14].** QSM is an online RL method that trains diffusion policies by matching the diffusion loss  
 229 to action gradients. QSM aims to avoid instability of value propagation to the expressive policy by  
 230 incorporating losses to guide the denoising process.

231 **DSRL [18].** DSRL is an online RL method that adapts a frozen diffusion-based behavior cloning  
 232 policy by performing RL over the initial noise seed used for sampling. DSRL learns a clipped  
 233 Gaussian policy that predicts the initial noise.

234 **QAM [37].** QAM is an online RL method that trains diffusion policies with adjoint matching. QAM  
 235 aims to avoid instability of value propagation to the expressive policy by using action gradients to  
 236 form a step-wise objective function that is free from unstable backpropagation.

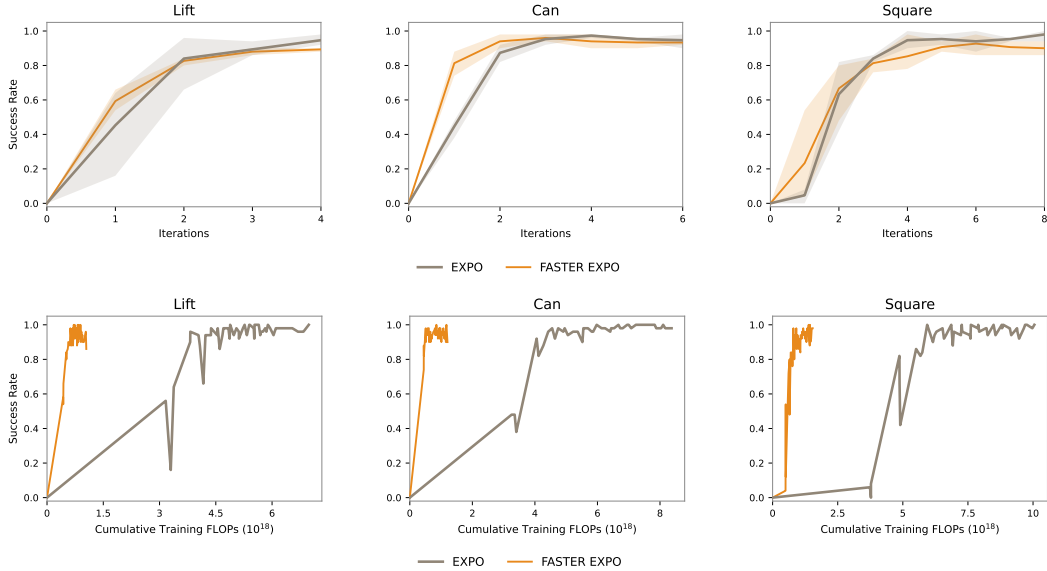


Figure 5: **Top: Success rate curves** of **FASTER-EXPO** and **EXPO** in the batch-online setting. **FASTER** matches the performance of **EXPO** in iterations. **Bottom: Compute comparisons** of **FASTER-EXPO** and **EXPO** in the batch-online setting. Like in the online setting, in the batch-online setting **FASTER-EXPO** yields a large FLOP reduction compared to **EXPO** from not needing to denoise all action samples.

237 **FQL [12]**. FQL uses a one-step flow policy to maximize the Q estimates learned with the standard  
 238 TD objective. It also incorporates a behavioral regularization term toward a BC flow policy.

239 **5.1 How does FASTER perform compared to state-of-the-art RL methods in both the online**  
 240 **setting and the batch-online setting?**

241 We present the main results in [Figure 3](#) and [Figure 5](#). Across both the online and batch-online  
 242 settings, **FASTER-EXPO** achieves the strongest overall performance among all baselines. In the online  
 243 setting, **FASTER-EXPO** significantly outperforms even highly sample-efficient methods such as RLPD.  
 244 Methods that rely on action gradients to directly train an expressive policy, such as QSM and QAM,  
 245 consistently struggle to match this performance. DSRL, which steers diffusion policies through  
 246 the noise latent space, improves steadily but exhibits worse sample efficiency than methods that  
 247 optimize directly in the action space. FQL, despite leveraging a one-step flow policy, consistently  
 248 underperforms as behavioral regularization toward a BC policy can hinder exploration in the online  
 249 setting. In the batch-online setting, **FASTER-EXPO** achieves the same performance as **EXPO** when  
 250 comparing iterations, and outperforms **EXPO** in terms of compute efficiency. Overall, **FASTER**  
 251 achieves state-of-the-art performance across all evaluated methods.

252 **5.2 Is FASTER able to recover the performance gains of test-time scaling while reducing the**  
 253 **computational cost?**

254 We now evaluate whether the noise-level critic can identify high-quality noise seeds *before* denoising,  
 255 avoiding the redundant computation of fully denoising all  $N$  candidates. We compare **FASTER-EXPO**  
 256 and **FASTER-IDQL** against their unfiltered counterparts, **EXPO** and **IDQL**, which rely on best-of- $N$   
 257 sampling at inference time. As shown in [Figure 3](#) and [Figure 4](#), both **FASTER** variants achieve task  
 258 success rates within the margin of error of their respective baselines across the evaluation suite in  
 259 the online setting, confirming that the noise-level critic effectively recovers the performance gains  
 260 of test-time scaling. Crucially, this comes at a fraction of the computational cost: as detailed in  
 261 [Section 4.4](#) and shown empirically in [Figure 3](#), [Figure 4](#) and [Figure 6](#), **FASTER** reduces inference-time  
 262 actor forward passes from  $\mathcal{O}(TN)$  to  $\mathcal{O}(T)$  by replacing  $N$  full denoising rollouts with  $N$  lightweight  
 263 critic evaluations followed by a single denoising pass. Concretely, the update step goes from 11.6 s to  
 264 2.5 s and the inference step goes from 566 ms to 335 ms. These results demonstrate that the sample  
 265 variance exploited by best-of- $N$  selection is largely determined at the noise level, and that **FASTER**  
 266 can capture this signal without incurring the cost of exhaustive denoising. We further discuss the  
 267 computational tradeoffs between our methods and baselines in [Section C.1](#).

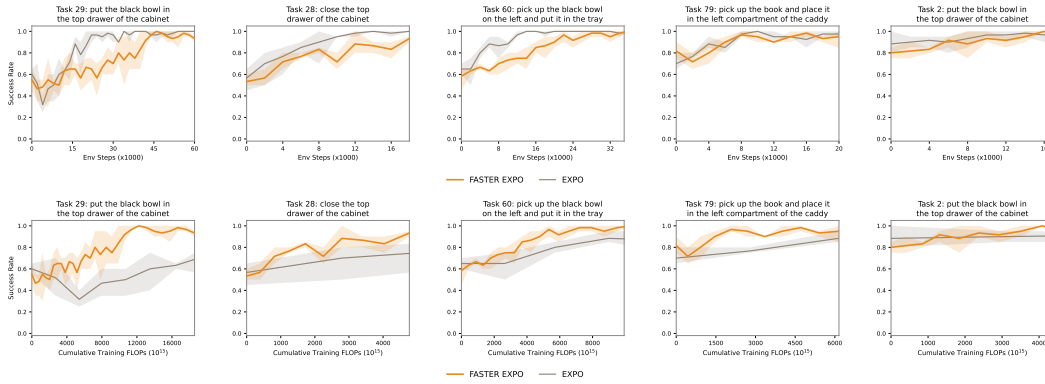


Figure 7: **FASTER-EXPO** compared to **EXPO** on top of  $\pi_{0.5}$ . **Top**: Performance with environment steps. **FASTER-EXPO** is competitive in performance compared to **EXPO**. **Bottom**: Performance with FLOPs. **FASTER-EXPO** performs significantly better than **EXPO** under the same compute as **FASTER-EXPO** chooses the best action sample without denoising all sampled actions in inference and training.

### 268 5.3 Can **FASTER** be applied to a pretrained Vision-Language-Action (VLA) model?

269 The computational burden of sampling-based  
 270 methods is particularly acute as model scale in-  
 271 creases; recent VLAs commonly reach 3B pa-  
 272 rameters, making the cost of denoising  $N$   
 273 candidates at every environment step prohibitive  
 274 for both training and deployment. We there-  
 275 fore evaluate whether **FASTER** can be applied  
 276 to a pretrained VLA to recover the benefits  
 277 of best-of- $N$  sampling at substantially lower  
 278 cost. Concretely, we apply **FASTER-EXPO**  
 279 to the 3.3B-parameter `pi05_libero` checkpoint  
 280 from OpenPI [30] and fine-tune it with online RL  
 281 on 5 held-out `libero_90` tasks, comparing against  
 282 **EXPO** under identical conditions.

282 **Compute savings.** As shown in Figure 6, **FASTER-EXPO** reduces the per-step update time by  
 283 approximately  $4.5\times$  relative to **EXPO**, from roughly 11.6 s to 2.5 s. This speedup stems from the fact  
 284 that **FASTER** only denoises a single candidate through the 3.3B-parameter actor during each  
 285 training step, whereas **EXPO** must denoise all  $N$  candidates to score them with the env-level critic.  
 286 At inference time, **FASTER-EXPO** reduces latency from 566 ms to 335 ms per action, a  $1.7\times$  speedup.  
 287 The gains are even more pronounced in FLOPs: **FASTER-EXPO** requires  $4.70 \times 10^{12}$  inference FLOPs  
 288 per step compared to  $3.75 \times 10^{13}$  for **EXPO**—an  $8\times$  reduction (Figure 6). The theoretical tradeoff  
 289 described in Section 4.4 is especially favorable in the VLA setting because the noise-level critic  
 290 (20M parameters) is far smaller than the actor (3.3B parameters), so scoring  $N$  noise candidates adds  
 291 negligible overhead relative to a single denoising pass.

292 **Task performance.** Despite these substantial compute reductions, **FASTER-EXPO** matches **EXPO** on  
 293 the majority of the 5 evaluated tasks (Figure 7, top). Across the task suite, **FASTER-EXPO** achieves  
 294 comparable aggregate success rates to **EXPO** while requiring a fraction of the compute budget. These  
 295 results demonstrate that **FASTER** scales naturally to large pretrained VLAs, enabling efficient online  
 296 RL fine-tuning without sacrificing the performance benefits of sampling-based action selection.

### 297 5.4 What components of **FASTER** are most important for performance?

298 To better understand how different pieces of **FASTER** contribute to performance, we ablate over three  
 299 key components: (1) the size of the denoise critic, (2) which step to filter with the denoise critic,  
 300 and (3) learning the full filtering policy according to the MDP in Section 4.1. We present additional  
 301 experiments on comparing **FASTER** against distilling the value maximizing distribution in Section A  
 302 as another approach for reducing computational cost.

303 **How does the size of the denoise critic affect performance?** We investigate whether the capacity  
 304 of the denoise critic is a significant factor in performance. Specifically, we compare architectures  
 305 ranging from a smaller network (128 hidden units  $\times$  3 layers, 300K params) to a medium network

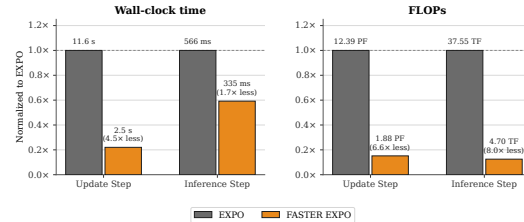


Figure 6: **Training and inference timing of FASTER-EXPO** compared to **EXPO**. **FASTER-EXPO** achieves  $1.7\times$  improvement in inference time and  $4.5\times$  improvement in the update step time.

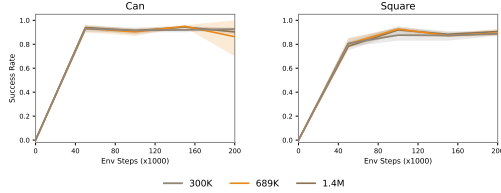


Figure 8: Critic-size ablation for **FASTER-EXPO** on can and square. We compare filtering critics  $Q^{dn}$  with parameter counts set to approximately  $1.0\times$ ,  $0.5\times$ , and  $0.25\times$  that of  $Q^a$ . We find that  $Q^{dn}$  can be substantially smaller than  $Q^a$  without degrading performance.

306 (176 hidden units  $\times$  3 layers, 689K params) to a larger network (256 hidden units  $\times$  3 layers, 1.4M  
 307 params) on the Robomimic Can and Square tasks. As shown in Figure 8, performance remains  
 308 largely consistent across all configurations, suggesting that **FASTER** is robust to moderate variations in  
 309 critic capacity and does not require careful tuning of this hyperparameter. Importantly, this suggests  
 310 **FASTER** can be used with a smaller denoise critic compared to action critic, which can greatly reduce  
 311 computational cost in the case of large critic networks.

312 **How does the choice of filtering step affect performance?** In our experiments, we im-  
 313 plement **FASTER** by filtering only at the initial noise. A natural question is whether the specific  
 314 denoising timestep at which the critic filter is applied has a meaningful impact on performance.  
 315 We evaluate several choices of filtering step on the Robomimic Can and Square tasks. As re-  
 316 ported in Figure 9, performance is consistent across the range of steps considered, indicating  
 317 that **FASTER** does not require precise selection of the filtering timestep and is robust to this design choice;  
 318 filtering at the initial seed achieves the best performance relative to the computation required.

325 **Policy learning with the full filtering MDP.** The complete formulation of **FASTER** learns a filtering  
 326 policy that can filter at any step of the denoising process. In our experiments, we reduce this multi-step  
 327 filtering process to a single decision and implement filtering at the initial noise level, as this approach  
 328 incurs minimal computational overhead while achieving performance on par with fully denoising  
 329 all action candidates. Here, we present results for the general setting in which the filtering policy  
 330 is learned to adaptively select the denoising step at which filtering occurs. In this case,  $\gamma$  implicitly  
 331 encourages the policy to terminate at earlier timesteps to save compute. Evaluating on the Robomimic  
 332 Can and Square tasks, we find that the learned adaptive policy performs comparably to fixed filtering  
 333 at the initial noise level. This suggests the single-step simplification of filtering at the initial seed  
 334 serves as a good proxy for solving the full MDP.

## 335 6 Discussion

336 In this work, we propose **FASTER**, a method for obtaining the benefits of sampling-based test-time  
 337 scaling without incurring its computational cost. By framing the denoising of multiple action candi-  
 338 dates as an MDP, **FASTER** learns a filtering policy that selects actions according to their downstream  
 339 value. Despite these promising results, **FASTER** has limitations. First, while **FASTER** substantially  
 340 improves upon its base algorithms in computational efficiency, it is not designed to improve sample  
 341 efficiency. In our experiments, **FASTER-EXPO** achieves better sample efficiency relative to strong  
 342 baselines, owing largely to the inherent sample efficiency of EXPO, but does not surpass EXPO itself  
 343 in this regard. Improving the sample efficiency of filtering-based approaches is an interesting direction  
 344 for future work. Furthermore, **FASTER** is designed for policy classes that use initial noise seeds,  
 345 which encompass some of the most performant policy representations in contemporary reinforcement  
 346 learning. Extending this to policy classes that lack such a seed structure is left to future work.

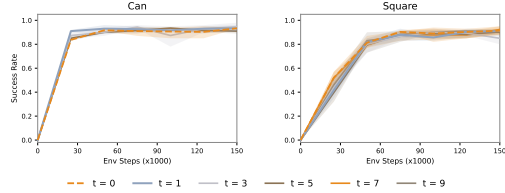


Figure 9: Filtering step ablation. **FASTER-EXPO** achieves similar performance across different filtering steps, suggesting **FASTER** does not require precise selection of the filtering step and filtering at the initial seed is effective.

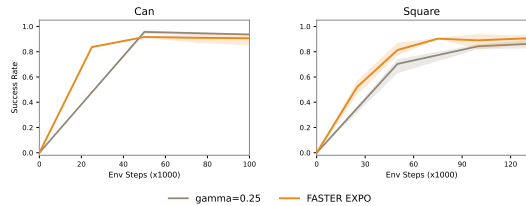


Figure 10: Learned-filter ablation results. Filtering at the initial seed performs comparably to learning the full filtering policy in the MDP.

## 347 References

- 348 [1] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha  
349 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language  
350 models. *CoRR*, abs/2203.11171, 2022. doi: 10.48550/arXiv.2203.11171. URL [https://](https://arxiv.org/abs/2203.11171)  
351 [arxiv.org/abs/2203.11171](https://arxiv.org/abs/2203.11171).
- 352 [2] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute  
353 optimally can be more effective than scaling model parameters. *CoRR*, abs/2408.03314, 2024.  
354 doi: 10.48550/arXiv.2408.03314. URL <https://arxiv.org/abs/2408.03314>.
- 355 [3] Yinlam Chow, Guy Tennenholtz, Izzeddin Gur, Vincent Zhuang, Bo Dai, Sridhar Thiagarajan,  
356 Craig Boutilier, Rishabh Agarwal, Aviral Kumar, and Aleksandra Faust. Inference-aware  
357 fine-tuning for best-of-N sampling in large language models. *CoRR*, abs/2412.15287, 2024. doi:  
358 10.48550/arXiv.2412.15287. URL <https://arxiv.org/abs/2412.15287>.
- 359 [4] Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan  
360 Yang, Yandong Li, Tommi Jaakkola, Xuhui Jia, and Saining Xie. Inference-time scaling  
361 for diffusion models beyond scaling denoising steps. *CoRR*, abs/2501.09732, 2025. doi:  
362 10.48550/arXiv.2501.09732. URL <https://arxiv.org/abs/2501.09732>.
- 363 [5] Perry Dong, Qiyang Li, Dorsa Sadigh, and Chelsea Finn. EXPO: Stable reinforcement learning  
364 with expressive policies. *CoRR*, abs/2507.07986, 2025. doi: 10.48550/arXiv.2507.07986. URL  
365 <https://arxiv.org/abs/2507.07986>.
- 366 [6] Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey  
367 Levine. IDQL: Implicit q-learning as an actor-critic method with diffusion policies. *CoRR*,  
368 abs/2304.10573, 2023. doi: 10.48550/arXiv.2304.10573. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2304.10573)  
369 [2304.10573](https://arxiv.org/abs/2304.10573).
- 370 [7] Perry Dong, Chongyi Zheng, Chelsea Finn, Dorsa Sadigh, and Benjamin Eysenbach. Value  
371 flows, 2026. URL <https://arxiv.org/abs/2510.07650>.
- 372 [8] Huayu Chen, Cheng Lu, Chengyang Ying, Hang Su, and Jun Zhu. Offline reinforcement  
373 learning via high-fidelity generative behavior modeling, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2209.14548)  
374 [abs/2209.14548](https://arxiv.org/abs/2209.14548).
- 375 [9] Zhendong Wang, Jonathan J. Hunt, and Mingyuan Zhou. Diffusion policies as an expressive  
376 policy class for offline reinforcement learning. *CoRR*, abs/2208.06193, 2022. doi: 10.48550/  
377 arXiv.2208.06193. URL <https://arxiv.org/abs/2208.06193>.
- 378 [10] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and  
379 Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *CoRR*,  
380 abs/2303.04137, 2023. doi: 10.48550/arXiv.2303.04137. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2303.04137)  
381 [2303.04137](https://arxiv.org/abs/2303.04137).
- 382 [11] Long Yang, Zhixiong Huang, Fenghao Lei, Yucun Zhong, Yiming Yang, Cong Fang, Shiting  
383 Wen, Binbin Zhou, and Zhouchen Lin. Policy representation via diffusion probability model for  
384 reinforcement learning, 2023. URL <https://arxiv.org/abs/2305.13122>.
- 385 [12] Seohong Park, Qiyang Li, and Sergey Levine. Flow q-learning. In Aarti Singh, Maryam Fazel,  
386 Daniel Hsu, Simon Lacoste-Julien, Felix Berkenkamp, Tegan Maharaj, Kiri Wagstaff, and  
387 Jerry Zhu, editors, *Forty-second International Conference on Machine Learning, ICML 2025,*  
388 *Vancouver, BC, Canada, July 13-19, 2025*, Proceedings of Machine Learning Research. PMLR /  
389 OpenReview.net, 2025. URL <https://proceedings.mlr.press/v267/park25f.html>.
- 390 [13] Andrew Wagenmaker, Perry Dong, Raymond Tsao, Chelsea Finn, and Sergey Levine. Posterior  
391 behavioral cloning: Pretraining bc policies for efficient rl finetuning, 2025. URL [https://](https://arxiv.org/abs/2512.16911)  
392 [arxiv.org/abs/2512.16911](https://arxiv.org/abs/2512.16911).
- 393 [14] Michael Psenka, Alejandro Escontrela, Pieter Abbeel, and Yi Ma. Learning a diffusion model  
394 policy from rewards via q-score matching. In *Proceedings of the 41st International Conference*  
395 *on Machine Learning*, 2024.

- 396 [15] Mitsuhiko Nakamoto, Oier Mees, Aviral Kumar, and Sergey Levine. Steering your generalists:  
397 Improving robotic foundation models via value guidance. *CoRR*, abs/2410.13816, 2024. doi:  
398 10.48550/arXiv.2410.13816. URL <https://arxiv.org/abs/2410.13816>.
- 399 [16] Jacky Kwok, Christopher Agia, Rohan Sinha, Matthew Foutter, Shulu Li, Ion Stoica, Azalia  
400 Mirhoseini, and Marco Pavone. Robomoney: Scaling test-time sampling and verification for  
401 vision-language-action models. *CoRR*, abs/2506.17811, 2025. doi: 10.48550/arXiv.2506.17811.  
402 URL <https://arxiv.org/abs/2506.17811>.
- 403 [17] Zhendong Wang, Zhaoshuo Li, Ajay Mandlekar, Zhenjia Xu, Jiaojiao Fan, Yashraj Narang,  
404 Linxi Fan, Yuke Zhu, Yogesh Balaji, Mingyuan Zhou, Ming-Yu Liu, and Yu Zeng. One-step  
405 diffusion policy: Fast visuomotor policies via diffusion distillation. *CoRR*, abs/2410.21257,  
406 2024. doi: 10.48550/arXiv.2410.21257. URL <https://arxiv.org/abs/2410.21257>.
- 407 [18] Andrew Wagenmaker, Mitsuhiko Nakamoto, Yunchu Zhang, Seohong Park, Waleed Yagoub,  
408 Anusha Nagabandi, Abhishek Gupta, and Sergey Levine. Steering your diffusion policy with  
409 latent space reinforcement learning. *CoRR*, abs/2506.15799, 2025. doi: 10.48550/arXiv.2506.  
410 15799. URL <https://arxiv.org/abs/2506.15799>.
- 411 [19] Yuanhao Ban, Ruochen Wang, Tianyi Zhou, Boqing Gong, Cho-Jui Hsieh, and Minhao Cheng.  
412 The crystal ball hypothesis in diffusion models: Anticipating object positions from initial noise.  
413 *CoRR*, abs/2406.01970, 2024. doi: 10.48550/arXiv.2406.01970. URL <https://arxiv.org/abs/2406.01970>.  
414 <https://arxiv.org/abs/2406.01970>.
- 415 [20] Zipeng Qi, Lichen Bai, Haoyi Xiong, and Zeke Xie. Not all noises are created equally: Diffusion  
416 noise selection and optimization. *CoRR*, abs/2407.14041, 2024. doi: 10.48550/arXiv.2407.  
417 14041. URL <https://arxiv.org/abs/2407.14041>.
- 418 [21] Xiefan Guo, Jinlin Liu, Miaomiao Cui, Jiankai Li, Hongyu Yang, and Di Huang. InitNO:  
419 Boosting text-to-image diffusion models via initial noise optimization. *CoRR*, abs/2404.04650,  
420 2024. doi: 10.48550/arXiv.2404.04650. URL <https://arxiv.org/abs/2404.04650>.
- 421 [22] Changgu Chen, Libing Yang, Xiaoyan Yang, Lianggangxu Chen, Gaoqi He, Changbo Wang,  
422 and Yang Li. FIND: Fine-tuning initial noise distribution with policy optimization for diffusion  
423 models. *CoRR*, abs/2407.19453, 2024. doi: 10.48550/arXiv.2407.19453. URL <https://arxiv.org/abs/2407.19453>.  
424 <https://arxiv.org/abs/2407.19453>.
- 425 [23] Zikai Zhou, Shitong Shao, Lichen Bai, Zhiqiang Xu, Bo Han, and Zeke Xie. Golden noise for  
426 diffusion models: A learning framework. *CoRR*, abs/2411.09502, 2024. doi: 10.48550/arXiv.  
427 2411.09502. URL <https://arxiv.org/abs/2411.09502>.
- 428 [24] Donghoon Ahn, Jiwon Kang, Sanghyun Lee, Jaewon Min, Minjae Kim, Wooseok Jang,  
429 Hyoungwon Cho, Sayak Paul, SeonHwa Kim, Eunju Cha, Kyong Hwan Jin, and Seun-  
430 gryong Kim. A noise is worth diffusion guidance. *CoRR*, abs/2412.03895, 2024. doi:  
431 10.48550/arXiv.2412.03895. URL <https://arxiv.org/abs/2412.03895>.
- 432 [25] Harvey Mannering, Zhiwu Huang, and Adam Prügel-Bennett. Noise-level diffusion guidance:  
433 Well begun is half done. *CoRR*, abs/2509.13936, 2025. doi: 10.48550/arXiv.2509.13936. URL  
434 <https://arxiv.org/abs/2509.13936>.
- 435 [26] Luca Eyring, Shyamgopal Karthik, Alexey Dosovitskiy, Nataniel Ruiz, and Zeynep Akata. Noise  
436 hypernetworks: Amortizing test-time compute in diffusion models. *CoRR*, abs/2508.09968,  
437 2025. doi: 10.48550/arXiv.2508.09968. URL <https://arxiv.org/abs/2508.09968>.
- 438 [27] Qingtao Yu, Changlin Song, Minghao Sun, Zhengyang Yu, Vinay Kumar Verma, Soumya Roy,  
439 Sumit Negi, Hongdong Li, and Dylan Campbell. TTSnap: Test-time scaling of diffusion models  
440 via noise-aware pruning. *CoRR*, abs/2511.22242, 2025. doi: 10.48550/arXiv.2511.22242. URL  
441 <https://arxiv.org/abs/2511.22242>.
- 442 [28] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate  
443 and transfer data with rectified flow, 2022. URL <https://arxiv.org/abs/2209.03003>.
- 444 [29] Kevin Frans, Danijar Hafner, Sergey Levine, and Pieter Abbeel. One step diffusion via shortcut  
445 models, 2025. URL <https://arxiv.org/abs/2410.12557>.

- 446 [30] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny  
447 Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, et al.  $\pi_{0,5}$ : A Vision-  
448 Language-Action Model with Open-World Generalization. *arXiv preprint arXiv:2504.16054*,  
449 2025.
- 450 [31] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini,  
451 Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion  
452 English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified  
453 flow transformers for high-resolution image synthesis, 2024. URL [https://arxiv.org/abs/  
454 2403.03206](https://arxiv.org/abs/2403.03206).
- 455 [32] Yu Gao, Lixue Gong, Qiushan Guo, Xiaoxia Hou, Zhichao Lai, Fanshi Li, Liang Li, Xiaochen  
456 Lian, Chao Liao, Liyang Liu, Wei Liu, Yichun Shi, Shiqi Sun, Yu Tian, Zhi Tian, Peng Wang,  
457 Rui Wang, Xuanda Wang, Xun Wang, Ye Wang, Guofeng Wu, Jie Wu, Xin Xia, Xuefeng Xiao,  
458 Zhonghua Zhai, Xinyu Zhang, Qi Zhang, Yuwei Zhang, Shijia Zhao, Jianchao Yang, and Weilin  
459 Huang. Seedream 3.0 technical report, 2025. URL <https://arxiv.org/abs/2504.11346>.
- 460 [33] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu,  
461 Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou,  
462 Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng,  
463 Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun,  
464 Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei  
465 Wang, Wenmeng Zhou, Wenteng Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming  
466 Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang,  
467 Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi,  
468 Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced  
469 large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.
- 470 [34] Physical Intelligence, Ali Amin, Raichelle Aniceto, Ashwin Balakrishna, Kevin Black, Ken  
471 Conley, Grace Connors, James Darpinian, Karan Dhabalia, Jared DiCarlo, Danny Driess,  
472 Michael Equi, Adnan Esmail, Yunhao Fang, Chelsea Finn, Catherine Glossop, Thomas Godden,  
473 Ivan Goryachev, Lachy Groom, Hunter Hancock, Karol Hausman, Gashon Hussein, Brian  
474 Ichter, Szymon Jakubczak, Rowan Jen, Tim Jones, Ben Katz, Liyiming Ke, Chandra Kuchi,  
475 Marinda Lamb, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Yao Lu, Vishnu Mano, Mohith  
476 Mothukuri, Suraj Nair, Karl Pertsch, Allen Z. Ren, Charvi Sharma, Lucy Xiaoyang Shi, Laura  
477 Smith, Jost Tobias Springenberg, Kyle Stachowicz, Will Stoeckle, Alex Swerdlow, James  
478 Tanner, Marcel Torne, Quan Vuong, Anna Walling, Haohuan Wang, Blake Williams, Sukwon  
479 Yoo, Lili Yu, Ury Zhilinsky, and Zhiyuan Zhou.  $\pi_{0,6}^*$ : a vla that learns from experience, 2025.  
480 URL <https://arxiv.org/abs/2511.14759>.
- 481 [35] Perry Dong, Suvir Mirchandani, Dorsa Sadigh, and Chelsea Finn. What matters for batch online  
482 reinforcement learning in robotics?, 2025. URL <https://arxiv.org/abs/2505.08078>.
- 483 [36] Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement  
484 learning with offline data. In *International Conference on Machine Learning*, pages 1577–1594.  
485 PMLR, 2023.
- 486 [37] Qiyang Li and Sergey Levine. Q-learning with adjoint matching, 2026. URL [https://arxiv.  
487 org/abs/2601.14234](https://arxiv.org/abs/2601.14234).
- 488 [38] Perry Dong, Alec M. Lessing, Annie S. Chen, and Chelsea Finn. Reinforcement learning via  
489 implicit imitation guidance, 2025. URL <https://arxiv.org/abs/2506.07505>.
- 490 [39] Perry Dong, Kuo-Han Hung, Alexander Swerdlow, Dorsa Sadigh, and Chelsea Finn. Tql:  
491 Scaling q-functions with transformers by preventing attention collapse, 2026. URL [https://  
492 arxiv.org/abs/2602.01439](https://arxiv.org/abs/2602.01439).
- 493 [40] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo  
494 Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming  
495 Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang  
496 Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky.  $\pi_0$ : A  
497 vision-language-action flow model for general robot control, 2026. URL [https://arxiv.  
498 org/abs/2410.24164](https://arxiv.org/abs/2410.24164).

- 499 [41] NVIDIA, Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding,  
500 Linxi "Jim" Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, Joel Jang, Zhenyu Jiang,  
501 Jan Kautz, Kaushil Kundalia, Lawrence Lao, Zhiqi Li, Zongyu Lin, Kevin Lin, Guilin Liu,  
502 Edith Llontop, Loic Magne, Ajay Mandlekar, Avnish Narayan, Soroush Nasiriany, Scott Reed,  
503 You Liang Tan, Guanzhi Wang, Zu Wang, Jing Wang, Qi Wang, Jiannan Xiang, Yuqi Xie,  
504 Yinzhen Xu, Zhenjia Xu, Seonghyeon Ye, Zhiding Yu, Ao Zhang, Hao Zhang, Yizhou Zhao,  
505 Ruijie Zheng, and Yuke Zhu. Gr00t n1: An open foundation model for generalist humanoid  
506 robots, 2025. URL <https://arxiv.org/abs/2503.14734>.
- 507 [42] Seonghyeon Ye, Yunhao Ge, Kaiyuan Zheng, Shenyuan Gao, Sihyun Yu, George Kurian, Suneel  
508 Indupuru, You Liang Tan, Chuning Zhu, Jiannan Xiang, Ayaan Malik, Kyungmin Lee, William  
509 Liang, Nadun Ranawaka, Jiasheng Gu, Yinzhen Xu, Guanzhi Wang, Fengyuan Hu, Avnish  
510 Narayan, Johan Bjorck, Jing Wang, Gwanghyun Kim, Dantong Niu, Ruijie Zheng, Yuqi Xie,  
511 Jimmy Wu, Qi Wang, Ryan Julian, Danfei Xu, Yilun Du, Yevgen Chebotar, Scott Reed, Jan  
512 Kautz, Yuke Zhu, Linxi "Jim" Fan, and Joel Jang. World action models are zero-shot policies,  
513 2026. URL <https://arxiv.org/abs/2602.15922>.
- 514 [43] Jonas Pai, Liam Achenbach, Victoriano Montesinos, Benedek Forrai, Oier Mees, and Elvis  
515 Nava. mimic-video: Video-action models for generalizable robot control beyond vlas, 2025.  
516 URL <https://arxiv.org/abs/2512.15692>.

517 **A Additional Experiments**

518 A natural question is whether the best-of- $N$  sampling benefits observed in EXPO and IDQL can  
 519 be recovered by a single policy trained via distillation. Concretely, an alternative to **FASTER** would  
 520 proceed as follows: sample  $N$  noise candidates  $\{\epsilon_i\}_{i=1}^N$  with  $\epsilon_i \sim \mathcal{N}(0, I)$ , denoise each to obtain  
 521 actions  $a_i = \pi_\theta(s, \epsilon_i)$ , select the highest-scoring index  $i^* = \arg \max_i Q^a(s, a_i)$ , and distill the  
 522 selected action  $a_{i^*}$  into a new policy  $\pi'_\theta$  by supervising against  $a_{i^*}$ . We evaluate this distillation  
 523 baseline in Figure 11 and find that it performs substantially worse than **FASTER**. We attribute this gap  
 524 to the instability inherent in distillation: the distilled policy must continually chase a moving target  
 525 defined by the  $Q$ -function, yielding a non-stationary training distribution that destabilizes learning.  
 526 By contrast, **FASTER** operates directly in noise space, framing selection as an easier filtering problem  
 527 that is more amenable to stable optimization.

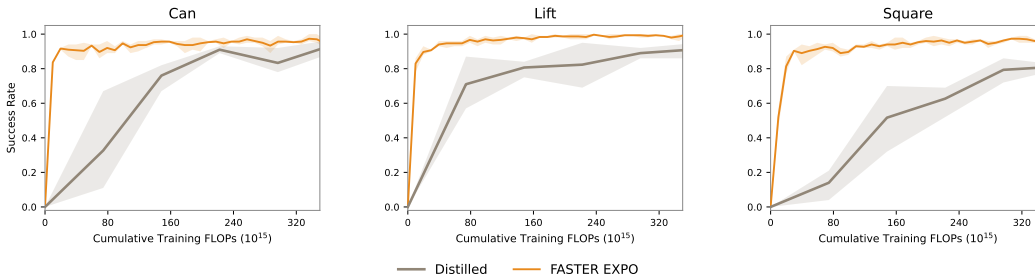


Figure 11: **Distillation ablation results.** Distilling the value maximizing distribution into a policy performs significantly worse compared to **FASTER**.

528 **B Experiment Details**

Hyperparameter	<b>FASTER-EXPO</b>	<b>FASTER-IDQL</b>
Actor / critic learning rate	$3 \times 10^{-4} / 3 \times 10^{-4}$	$3 \times 10^{-4} / 3 \times 10^{-4}$
Batch size	256	256
Discount	0.99	0.99
Target update $\tau$	0.005	0.005
UTD ratio	20	20
Candidates $N$	8	8
Candidates kept after filtering	1	1
Denosing steps $T$	10	100
Filter temperature mode	z-score	z-score
Expectile	—	0.8
Edit scale	0.15	—

Table 1: **Hyperparameters** for the **FASTER-EXPO** and **FASTER-IDQL** variants on Robomimic.

529 **Datasets and evaluation.** For the Robomimic experiments, all tasks use sparse binary rewards and  
 530 report success rate. In the online setting, we do not pretrain on offline data, but we retain the same task  
 531 conventions used in prior work to define task difficulty: **Lift** uses a 10-episode subset of the original  
 532 dataset to make it more challenging, **Can** uses the multi-human dataset, and **Square** and **Tool Hang**  
 533 use the standard proficient-human split. In the batch-online setting, we initialize training from these  
 534 same demonstration sources and then continue interaction online. For the batch-online setting, we use  
 535  $T = 10$  denosing steps for **Lift** and **Can**, and  $T = 100$  denosing steps for **Square**. For the final  
 536 rerun **FASTER-EXPO** experiments, evaluation is performed every 50k environment steps on **Lift**,  
 537 **Can**, and **Square**, and every 100k steps on **Tool Hang**; **FASTER-IDQL** uses the same evaluation  
 538 protocol. Evaluation-time filtering with  $Q^{dn}$  is greedy (temperature 0) for all Robomimic tasks,  
 539 while at training time we sample candidates from a softmax over z-score-normalized  $Q$ -values with  
 540 temperature 1.0 for **Lift**, **Can**, and **Square**, and 0.1 for **Tool Hang**. Unless otherwise specified,  
 541 line plots report the mean success rate over three seeds and the shaded regions denote the max and  
 542 min.

543 **VLA experiments.** For the pretrained VLA experiments, we follow the held-out LIBERO protocol  
544 described in the main text and fine-tune entirely online, without mixing in offline data. We use a  
545 batch size of 32, UTD ratio 10, and 100k online training steps. We use  $N = 8$  candidates, retain a  
546 single candidate after filtering for **FASTER-EXPO**, and denoise for  $T = 10$  steps. The VLA predicts  
547 action chunks over 10 steps and we execute the entire chunk before replanning. Online optimization  
548 is performed at the episode level rather than every environment step: after each completed episode,  
549 we run three updates on newly sampled replay minibatches, with each update containing UTD 10  
550 critic updates. To best utilize the replay data, we train on all sliding window chunks. For this  
551 setting, the critic is lightweight relative to the actor, so the additional cost of evaluating multiple noise  
552 candidates remains small compared to a full multi-candidate denoising pass. Both training-time and  
553 evaluation-time candidate selection use greedy sampling (temperature 0), and we use an edit scale of  
554 0.2.

555 **Implementation Details.** We run experiments on a mix of NVIDIA A40, L40S, and H100 GPUs.  
556 All timing experiments are conducted on L40S GPUs. For update step timing, we consider the  
557 time taken for a single critic update, with other parameters as described in the prior paragraph. For  
558 inference step timing, we report BS=1 inference.

## 559 C Computational Profile

### 560 C.1 Comparison to non-diffusion-based methods

561 Our work focuses on capturing the benefits of sampling-based methods for diffusion models without  
562 incurring the costs of these approaches. Diffusion and flow-based models are the dominant policy  
563 class in robotics, particularly in the real-world with VLAs [40, 41] and WAMs [42, 43], hence our  
564 focus. Other methods, such as the Gaussian policies used in RLPD [36], are generally more efficient  
565 due to lower typical parameter counts and no iterative computation—i.e., only a single forward step at  
566 inference. However, their limited expressivity makes them ill-suited for many real-world applications.  
567 Thus, this work focuses on the performance of diffusion methods with the same parameter count,  
568 which we believe to be the most apt comparison.

## 569 **NeurIPS Paper Checklist**

### 570 • **Claims**

571 Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s  
572 contributions and scope?

573 Answer: [Yes]

574 Justification: The claims we make are supported by the empirical results. We demonstrate across  
575 multiple algorithms and settings that our method achieves better compute efficiency compared to  
576 baselines.

577 Guidelines:

- 578 • The answer [N/A] means that the abstract and introduction do not include the claims made in the  
579 paper.
- 580 • The abstract and/or introduction should clearly state the claims made, including the contributions  
581 made in the paper and important assumptions and limitations. A [No] or [N/A] answer to this  
582 question will not be perceived well by the reviewers.
- 583 • The claims made should match theoretical and experimental results, and reflect how much the  
584 results can be expected to generalize to other settings.
- 585 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not  
586 attained by the paper.

### 587 • **Limitations**

588 Question: Does the paper discuss the limitations of the work performed by the authors?

589 Answer: [Yes]

590 Justification: The paper discusses the limitations of the method which include not improving upon  
591 sample efficiency and of being designed for policy classes that use initial seeds such as diffusion or  
592 flow-based policies.

593 Guidelines:

- 594 • The answer [N/A] means that the paper has no limitation while the answer [No] means that the  
595 paper has limitations, but those are not discussed in the paper.
- 596 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 597 • The paper should point out any strong assumptions and how robust the results are to violations of  
598 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,  
599 asymptotic approximations only holding locally). The authors should reflect on how these  
600 assumptions might be violated in practice and what the implications would be.
- 601 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested  
602 on a few datasets or with a few runs. In general, empirical results often depend on implicit  
603 assumptions, which should be articulated.
- 604 • The authors should reflect on the factors that influence the performance of the approach. For  
605 example, a facial recognition algorithm may perform poorly when image resolution is low or  
606 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide  
607 closed captions for online lectures because it fails to handle technical jargon.
- 608 • The authors should discuss the computational efficiency of the proposed algorithms and how they  
609 scale with dataset size.
- 610 • If applicable, the authors should discuss possible limitations of their approach to address problems  
611 of privacy and fairness.
- 612 • While the authors might fear that complete honesty about limitations might be used by reviewers  
613 as grounds for rejection, a worse outcome might be that reviewers discover limitations that  
614 aren’t acknowledged in the paper. The authors should use their best judgment and recognize  
615 that individual actions in favor of transparency play an important role in developing norms that  
616 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize  
617 honesty concerning limitations.

### 618 • **Theory assumptions and proofs**

619 Question: For each theoretical result, does the paper provide the full set of assumptions and a  
620 complete (and correct) proof?

621 Answer: [N/A]

622 Justification: The paper does not include theoretical results.

623 Guidelines:

- 624 • The answer [N/A] means that the paper does not include theoretical results.

- 625 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 626 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 627 • The proofs can either appear in the main paper or the supplemental material, but if they appear in
- 628 the supplemental material, the authors are encouraged to provide a short proof sketch to provide
- 629 intuition.
- 630 • Inversely, any informal proof provided in the core of the paper should be complemented by
- 631 formal proofs provided in appendix or supplemental material.
- 632 • Theorems and Lemmas that the proof relies upon should be properly referenced.

633 • **Experimental result reproducibility**

634 Question: Does the paper fully disclose all the information needed to reproduce the main experi-  
635 mental results of the paper to the extent that it affects the main claims and/or conclusions of the  
636 paper (regardless of whether the code and data are provided or not)?

637 Answer: [Yes]

638 Justification: We describe the algorithm and practical implementation of our method fully, provide  
639 detailed hyperparameters in the appendix, and include an anonymous repository.

640 Guidelines:

- 641 • The answer [N/A] means that the paper does not include experiments.
- 642 • If the paper includes experiments, a [No] answer to this question will not be perceived well by
- 643 the reviewers: Making the paper reproducible is important, regardless of whether the code and
- 644 data are provided or not.
- 645 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make
- 646 their results reproducible or verifiable.
- 647 • Depending on the contribution, reproducibility can be accomplished in various ways. For
- 648 example, if the contribution is a novel architecture, describing the architecture fully might suffice,
- 649 or if the contribution is a specific model and empirical evaluation, it may be necessary to either
- 650 make it possible for others to replicate the model with the same dataset, or provide access to
- 651 the model. In general, releasing code and data is often one good way to accomplish this, but
- 652 reproducibility can also be provided via detailed instructions for how to replicate the results,
- 653 access to a hosted model (e.g., in the case of a large language model), releasing of a model
- 654 checkpoint, or other means that are appropriate to the research performed.
- 655 • While NeurIPS does not require releasing code, the conference does require all submissions
- 656 to provide some reasonable avenue for reproducibility, which may depend on the nature of the
- 657 contribution. For example
  - 658 • If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce
  - 659 that algorithm.
  - 660 • If the contribution is primarily a new model architecture, the paper should describe the
  - 661 architecture clearly and fully.
  - 662 • If the contribution is a new model (e.g., a large language model), then there should either be
  - 663 a way to access this model for reproducing the results or a way to reproduce the model (e.g.,
  - 664 with an open-source dataset or instructions for how to construct the dataset).
  - 665 • We recognize that reproducibility may be tricky in some cases, in which case authors are
  - 666 welcome to describe the particular way they provide for reproducibility. In the case of closed-
  - 667 source models, it may be that access to the model is limited in some way (e.g., to registered
  - 668 users), but it should be possible for other researchers to have some path to reproducing or
  - 669 verifying the results.

670 • **Open access to data and code**

671 Question: Does the paper provide open access to the data and code, with sufficient instructions to

672 faithfully reproduce the main experimental results, as described in supplemental material?

673 Answer: [Yes]

674 Justification: We provide an anonymous repository containing the code.

675 Guidelines:

- 676 • The answer [N/A] means that paper does not include experiments requiring code.
- 677 • Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.

678

- 679 • While we encourage the release of code and data, we understand that this might not be possible,  
680 so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless  
681 this is central to the contribution (e.g., for a new open-source benchmark).
- 682 • The instructions should contain the exact command and environment needed to run to reproduce  
683 the results. See the NeurIPS code and data submission guidelines ([https://neurips.cc/  
684 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 685 • The authors should provide instructions on data access and preparation, including how to access  
686 the raw data, preprocessed data, intermediate data, and generated data, etc.
- 687 • The authors should provide scripts to reproduce all experimental results for the new proposed  
688 method and baselines. If only a subset of experiments are reproducible, they should state which  
689 ones are omitted from the script and why.
- 690 • At submission time, to preserve anonymity, the authors should release anonymized versions (if  
691 applicable).
- 692 • Providing as much information as possible in supplemental material (appended to the paper) is  
693 recommended, but including URLs to data and code is permitted.
- 694 • **Experimental setting/details**  
695 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,  
696 how they were chosen, type of optimizer) necessary to understand the results?  
697 Answer: [Yes]  
698 Justification: We present all training and test details necessary to understand and replicate the  
699 results  
700 Guidelines:  
701 • The answer [N/A] means that the paper does not include experiments.  
702 • The experimental setting should be presented in the core of the paper to a level of detail that is  
703 necessary to appreciate the results and make sense of them.  
704 • The full details can be provided either with the code, in appendix, or as supplemental material.
- 705 • **Experiment statistical significance**  
706 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
707 information about the statistical significance of the experiments?  
708 Answer: [Yes]  
709 Justification: We include confidence intervals for plots in the paper.  
710 Guidelines:  
711 • The answer [N/A] means that the paper does not include experiments.  
712 • The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals,  
713 or statistical significance tests, at least for the experiments that support the main claims of the  
714 paper.  
715 • The factors of variability that the error bars are capturing should be clearly stated (for example,  
716 train/test split, initialization, random drawing of some parameter, or overall run with given  
717 experimental conditions).  
718 • The method for calculating the error bars should be explained (closed form formula, call to a  
719 library function, bootstrap, etc.)  
720 • The assumptions made should be given (e.g., Normally distributed errors).  
721 • It should be clear whether the error bar is the standard deviation or the standard error of the mean.  
722 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report  
723 a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is  
724 not verified.  
725 • For asymmetric distributions, the authors should be careful not to show in tables or figures  
726 symmetric error bars that would yield results that are out of range (e.g., negative error rates).  
727 • If error bars are reported in tables or plots, the authors should explain in the text how they were  
728 calculated and reference the corresponding figures or tables in the text.
- 729 • **Experiments compute resources**  
730 Question: For each experiment, does the paper provide sufficient information on the computer  
731 resources (type of compute workers, memory, time of execution) needed to reproduce the experi-  
732 ments?  
733 Answer: [Yes]  
734 Justification: We describe the compute resources used to run the experiments in the appendix.

735 Guidelines:

736 • The answer [N/A] means that the paper does not include experiments.

737 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud

738 provider, including relevant memory and storage.

739 • The paper should provide the amount of compute required for each of the individual experimental

740 runs as well as estimate the total compute.

741 • The paper should disclose whether the full research project required more compute than the

742 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into

743 the paper).

744 • **Code of ethics**

745 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS

746 Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

747 Answer: [Yes]

748 Justification: The research conducted in the paper conforms to the NeurIPS code of ethics.

749 Guidelines:

750 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.

751 • If the authors answer [No], they should explain the special circumstances that require a deviation

752 from the Code of Ethics.

753 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due

754 to laws or regulations in their jurisdiction).

755 • **Broader impacts**

756 Question: Does the paper discuss both potential positive societal impacts and negative societal

757 impacts of the work performed?

758 Answer: [N/A]

759 Justification: The work is a new algorithm for faster reinforcement learning on diffusion and flow-

760 matching models which does not pose specific societal impacts beyond improving the efficiency of

761 general reinforcement learning algorithms.

762 Guidelines:

763 • The answer [N/A] means that there is no societal impact of the work performed.

764 • If the authors answer [N/A] or [No], they should explain why their work has no societal impact

765 or why the paper does not address societal impact.

766 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,

767 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-

768 ment of technologies that could make decisions that unfairly impact specific groups), privacy

769 considerations, and security considerations.

770 • The conference expects that many papers will be foundational research and not tied to particular

771 applications, let alone deployments. However, if there is a direct path to any negative applications,

772 the authors should point it out. For example, it is legitimate to point out that an improvement in

773 the quality of generative models could be used to generate Deepfakes for disinformation. On the

774 other hand, it is not needed to point out that a generic algorithm for optimizing neural networks

775 could enable people to train models that generate Deepfakes faster.

776 • The authors should consider possible harms that could arise when the technology is being used

777 as intended and functioning correctly, harms that could arise when the technology is being used

778 as intended but gives incorrect results, and harms following from (intentional or unintentional)

779 misuse of the technology.

780 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies

781 (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitor-

782 ing misuse, mechanisms to monitor how a system learns from feedback over time, improving the

783 efficiency and accessibility of ML).

784 • **Safeguards**

785 Question: Does the paper describe safeguards that have been put in place for responsible release of

786 data or models that have a high risk for misuse (e.g., pre-trained language models, image generators,

787 or scraped datasets)?

788 Answer: [N/A]

789 Justification: The paper poses no such risks.

790 Guidelines:

- 791 • The answer [N/A] means that the paper poses no such risks.
- 792 • Released models that have a high risk for misuse or dual-use should be released with necessary
- 793 safeguards to allow for controlled use of the model, for example by requiring that users adhere to
- 794 usage guidelines or restrictions to access the model or implementing safety filters.
- 795 • Datasets that have been scraped from the Internet could pose safety risks. The authors should
- 796 describe how they avoided releasing unsafe images.
- 797 • We recognize that providing effective safeguards is challenging, and many papers do not require
- 798 this, but we encourage authors to take this into account and make a best faith effort.

799 • **Licenses for existing assets**

800 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,

801 properly credited and are the license and terms of use explicitly mentioned and properly respected?

802 Answer: [Yes]

803 Justification: We cite all papers, code, data, models we used, including in the code repository.

804 Guidelines:

- 805 • The answer [N/A] means that the paper does not use existing assets.
- 806 • The authors should cite the original paper that produced the code package or dataset.
- 807 • The authors should state which version of the asset is used and, if possible, include a URL.
- 808 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 809 • For scraped data from a particular source (e.g., website), the copyright and terms of service of
- 810 that source should be provided.
- 811 • If assets are released, the license, copyright information, and terms of use in the package should
- 812 be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for
- 813 some datasets. Their licensing guide can help determine the license of a dataset.
- 814 • For existing datasets that are re-packaged, both the original license and the license of the derived
- 815 asset (if it has changed) should be provided.
- 816 • If this information is not available online, the authors are encouraged to reach out to the asset's
- 817 creators.

818 • **New assets**

819 Question: Are new assets introduced in the paper well documented and is the documentation

820 provided alongside the assets?

821 Answer: [Yes]

822 Justification: New assets introduced in the paper are well documented

823 Guidelines:

- 824 • The answer [N/A] means that the paper does not release new assets.
- 825 • Researchers should communicate the details of the dataset/code/model as part of their submissions
- 826 via structured templates. This includes details about training, license, limitations, etc.
- 827 • The paper should discuss whether and how consent was obtained from people whose asset is
- 828 used.
- 829 • At submission time, remember to anonymize your assets (if applicable). You can either create an
- 830 anonymized URL or include an anonymized zip file.

831 • **Crowdsourcing and research with human subjects**

832 Question: For crowdsourcing experiments and research with human subjects, does the paper include

833 the full text of instructions given to participants and screenshots, if applicable, as well as details

834 about compensation (if any)?

835 Answer: [N/A]

836 Justification: The paper does not involve crowdsourcing or research with human subjects

837 Guidelines:

- 838 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with human
- 839 subjects.
- 840 • Including this information in the supplemental material is fine, but if the main contribution of the
- 841 paper involves human subjects, then as much detail as possible should be included in the main
- 842 paper.
- 843 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other
- 844 labor should be paid at least the minimum wage in the country of the data collector.

845 • **Institutional review board (IRB) approvals or equivalent for research with human subjects**

846 Question: Does the paper describe potential risks incurred by study participants, whether such  
847 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an  
848 equivalent approval/review based on the requirements of your country or institution) were obtained?

849 Answer: [N/A]

850 Justification: The paper does not involve crowdsourcing or research with human subjects.

851 Guidelines:

- 852 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with human  
853 subjects.
- 854 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be  
855 required for any human subjects research. If you obtained IRB approval, you should clearly state  
856 this in the paper.
- 857 • We recognize that the procedures for this may vary significantly between institutions and locations,  
858 and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their  
859 institution.
- 860 • For initial submissions, do not include any information that would break anonymity (if applicable),  
861 such as the institution conducting the review.

862 • **Declaration of LLM usage**

863 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard  
864 component of the core methods in this research? Note that if the LLM is used only for writing,  
865 editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or  
866 originality of the research, declaration is not required.

867 Answer: [N/A]

868 Justification: The core method development in this research does not involve LLMs as important,  
869 original, or non-standard components.

870 Guidelines:

- 871 • The answer [N/A] means that the core method development in this research does not involve  
872 LLMs as any important, original, or non-standard components.
- 873 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not be  
874 described.