
Can K Heads Explore Better Than One in Online Reinforcement Learning?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Ensemble methods provide principled exploration for Gaussian policies in rein-
2 forcement learning, but their extension to the increasingly popular family of gener-
3 ative policies based on diffusion and flow matching is non-trivial. The standard
4 ensemble template, namely K value heads trained on bootstrapped data subsets and
5 selected by exploration strategies, fails for Q-weighted generative losses. When the
6 K policy heads share a common Q-network, the Q-weighted importance weights
7 are identical across heads by construction, so the denoising loss produces head
8 invariant gradients in expectation and the heads collapse to a single function. Per-
9 sample bootstrap masking cannot prevent this collapse because the mask only
10 modulates a signal that is already head invariant before it enters the loss. This
11 failure mode is specific to importance weighted generative losses and does not arise
12 for standard Gaussian policy ensembles. To resolve this, we introduce Ensemble
13 Soft Actor-Critic (ESAC), an ensemble framework that pairs each generative head
14 with its own independent Q-network, so that per-head Q-weighted objectives drive
15 the heads toward distinct optima. Within the same architecture we propose Tra-
16 jectory Integrated Disagreement Exploration (TIDE), an action-selection rule that
17 scores candidates by ensemble disagreement integrated across the full denoising
18 trajectory and adaptively calibrates the bonus by the running correlation between
19 velocity disagreement and Q-uncertainty. We formalize the collapse failure mode
20 and the diversity restoration mechanism as theorems. On MuJoCo locomotion
21 benchmarks with both diffusion and conditional flow matching backbones, ESAC
22 delivers consistent multi-seed gains over its soft diffusion actor-critic backbone
23 with the largest gains on tasks where directed exploration matters most.

24 1 Introduction

25 The choice of policy representation fundamentally shapes what an RL agent can learn. Gaussian
26 policies, as used by SAC [Haarnoja et al., 2018] and its variants, parametrize actions as a diagonal
27 Gaussian with network-predicted mean and variance. This is inherently unimodal, the policy expresses
28 only one strategy per state, even when multiple qualitatively different actions are equally viable. The
29 limitation bites in manipulation, where distinct grasp strategies may all succeed, and in locomotion,
30 where multiple stable gaits coexist.

31 Generative policies based on diffusion [Ho et al., 2020, Song et al., 2021] and flow matching [Lipman
32 et al., 2023, Liu et al., 2023] close this expressiveness gap by replacing the parametric Gaussian
33 with an iterative sampling process, and recent work shows substantial empirical gains [Ma et al.,
34 2025, Ding et al., 2024, McAllister et al., 2025, Wang et al., 2023, Hansen-Estruch et al., 2024,
35 Ada et al., 2024]. Yet *exploration* in these methods remains primitive. Every published flow or
36 diffusion RL algorithm explores by adding isotropic Gaussian noise to the selected action or by
37 inheriting SAC’s entropy regularization, perturbing actions uniformly regardless of where the agent’s
38 uncertainty actually lies. On smooth reward landscapes like HalfCheetah this suffices, but on tasks
39 with coordinated multi-joint control and early termination, such as Ant-v4 or Humanoid-v4, it
40 is catastrophically sample inefficient. The Bayesian RL literature offers a principled alternative:
41 Bootstrapped DQN [Osband et al., 2016] maintains K value heads trained on bootstrapped experience
42 and selects one per episode via Thompson Sampling [Thompson, 1933], producing temporally
43 coherent exploration. SUNRISE [Lee et al., 2021] extended this to continuous control using ensemble
44 disagreement as a UCB-style [Auer et al., 2002] uncertainty proxy, and REDQ [Chen et al., 2021]
45 showed that Q-ensembles support stable high-UTD training. But these methods operate on Gaussian

46 policies [Fujimoto et al., 2018, Lillicrap et al., 2016], and their extension to generative policy classes
47 is non-trivial.

48 The naive extension fails. We show theoretically and empirically that any ensemble of generative
49 policy heads trained with a Q-weighted denoising loss against a *shared* Q-network collapses to a
50 single function regardless of bootstrap masking, independent network initialization, or architectural
51 separation. The mechanism is structural: the Q-weighted denoising gradient depends on the policy
52 heads only through the candidate actions they sample, and a shared Q-landscape produces head-
53 invariant gradients in expectation, so independent initialization shifts only the early dynamics rather
54 than the long-run optimum. Bootstrap masking adds per-sample noise but this variance is dwarfed
55 by the consistent mean signal from the shared Q, and Thompson Sampling degenerates to random
56 selection among K near-identical policies. The collapse is a property of the Q-weighted objective
57 itself and applies to any generative policy class trained by importance-weighted regression on a
58 shared critic.

59 *ESAC* resolves this by training a small team of K independent (policy, Q) pairs. Each head k optimizes
60 a Q-weighted objective using its *own* Q_{ϕ_k} for the importance weights, so different Q-landscapes drive
61 different policies even though all heads are trained on the same data. What independence buys is not
62 a different asymptote but a different *trajectory* through policy space, and the heads explore distinct
63 regions early and concentrate the team’s coverage of the action manifold. At every environment
64 step each policy proposes a small batch of candidate actions, the team scores every candidate by its
65 ensemble-mean return plus a disagreement bonus measured along the action-generation trajectory,
66 and the highest-scoring action is executed. The disagreement bonus, which we call *TIDE* (Trajectory
67 Integrated Disagreement Exploration), measures inter-head velocity-field disagreement at multiple
68 points along the denoising trajectory rather than only at the final action. Two heads can produce nearly
69 identical final actions via very different denoising paths; conversely, two heads can produce different
70 final actions by chance even when their underlying posteriors agree. Trajectory-level disagreement
71 separates these cases, and we calibrate the bonus strength online by the running Pearson correlation
72 between trajectory disagreement and value uncertainty so that the bonus is active only when the two
73 signals agree. The same recipe works with diffusion or flow-matching backbones with no further
74 change.

75 *ESAC* is the first ensemble exploration framework for expressive generative policies, and the first to
76 identify and resolve the shared-critic collapse mode. This combination unlocks directed, posterior-
77 driven exploration over multimodal action distributions in continuous control, with consistent gains on
78 hard locomotion tasks where undirected exploration plateaus and unimodal policies cannot represent
79 the multiple viable gaits. The key contributions for this work are as follows:

- 80 1. **Collapse theorem and architectural fix for generative-policy ensembles.** We prove that
81 any ensemble of generative heads trained with a Q-weighted loss against a shared critic
82 collapses to a single function regardless of initialization or bootstrap masking, and identify
83 per-head independent critics with circular twin-Q pessimism as the minimal architectural
84 change that breaks this collapse. The result generalizes to any generative policy class trained
85 by importance-weighted regression.
- 86 2. **Trajectory-integrated disagreement exploration (TIDE).** We introduce an action-
87 selection rule that exploits the iterative structure of generative sampling by measuring
88 inter-head disagreement at multiple denoising timesteps rather than only at the output.
89 Bonus strength is gated online by the correlation between disagreement and Q-uncertainty,
90 preventing over-exploration on unimodal landscapes while preserving gains where genuine
91 posterior diversity exists.
- 92 3. **Backbone comparison and modality-aware analysis.** We provide the systematic com-
93 parison of diffusion and conditional flow-matching backbones inside a unified ensemble
94 framework across five MuJoCo locomotion benchmarks and four additional tasks, showing
95 flow matching maintains substantially higher inter-head diversity throughout training, and
96 characterize when each backbone wins via Q-landscape modality analysis.

97 2 Preliminaries

98 **Entropy Regularized Reinforcement Learning.** We consider the standard infinite-horizon dis-
99 counted MDP $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$ with continuous state space $\mathcal{S} \subseteq \mathbb{R}^{d_s}$, continuous action space $\mathcal{A} \subseteq \mathbb{R}^{d_a}$,
100 transition dynamics $P(s'|s, a)$, reward function $r(s, a)$, and discount factor $\gamma \in (0, 1)$. Following
101 the maximum entropy RL framework [Haarnoja et al., 2018], we seek a policy π maximizing the
102 entropy-augmented objective

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \tau_{\text{ent}} \mathcal{H}(\pi(\cdot|s_t))) \right] \quad (1)$$

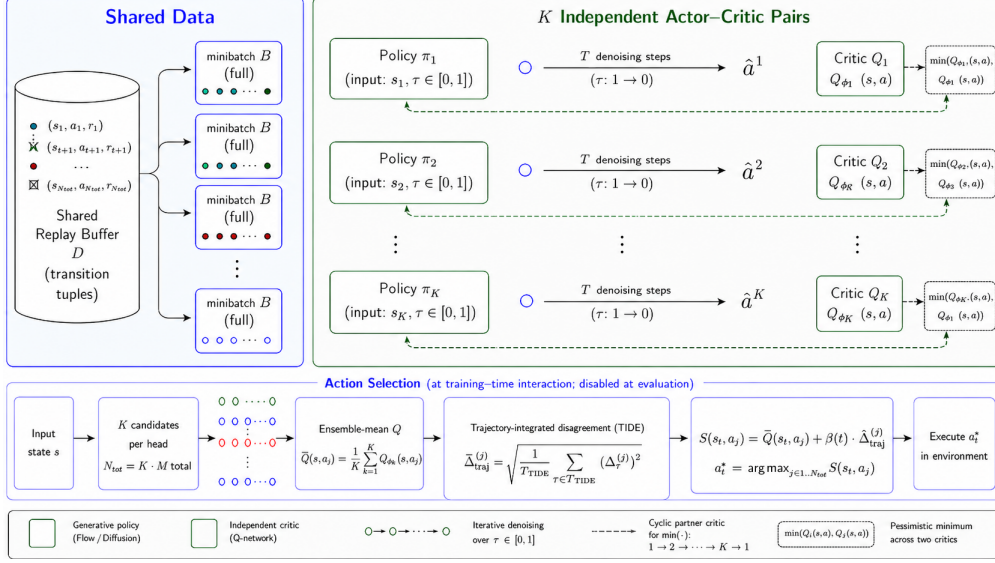


Figure 1: **ESAC training and inference pipeline.** A shared replay-buffer minibatch is fed to K independent actor-critic pairs. Each pair $(\pi_{\theta_k}, Q_{\phi_k})$ is trained jointly: the Q-loss uses a REDQ-style random subset minimum across the K Q-networks as a shared TD target, while the policy loss uses per-head importance weights computed via circular twin-Q pairing, so that each head’s policy update sees a different pair of Q-networks. At interaction time (bottom), every head proposes candidate actions and TIDE selects one by combining the ensemble Q-mean with a trajectory-integrated inter-head disagreement bonus.

103 where $\tau_{\text{ent}} > 0$ is a learned entropy temperature coefficient and $\mathcal{H}(\pi(\cdot|s)) = -\mathbb{E}_{a \sim \pi}[\log \pi(a|s)]$ is
 104 the policy entropy. SAC and its derivatives [Haarnoja et al., 2018, Fujimoto et al., 2018] form the
 105 foundation for modern continuous control. The soft Q-function satisfies the soft Bellman equation

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P} [V^\pi(s')], \quad V^\pi(s) = \mathbb{E}_{a \sim \pi} [Q^\pi(s, a) - \tau_{\text{ent}} \log \pi(a|s)] \quad (2)$$

106 In practice, SAC [Haarnoja et al., 2018] maintains two Q-networks Q_{ϕ_1}, Q_{ϕ_2} (twin-Q) and uses
 107 the minimum for both policy optimization and Bellman targets to suppress overestimation bias.
 108 The target Q-networks $Q_{\bar{\phi}_i}$ are updated via exponential moving average with coefficient ρ , namely
 109 $\bar{\phi}_i \leftarrow (1 - \rho)\bar{\phi}_i + \rho \phi_i$. For numerical stability, implementations parametrize $\log \tau_{\text{ent}}$ rather than τ_{ent}
 110 directly and recover the temperature via exponentiation.

111 **Diffusion based policies** SDAC [Ma et al., 2025] parametrizes the policy as a denoising diffusion
 112 probabilistic model (DDPM) [Ho et al., 2020]. Given state s , action candidates are generated by
 113 starting from Gaussian noise $x_T \sim \mathcal{N}(0, I)$ and iteratively applying a learned denoiser ϵ_θ over T
 114 diffusion steps,

$$x_{k-1} = \frac{1}{\sqrt{\alpha_k}} \left(x_k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k}} \epsilon_\theta(x_k, k, s) \right) + \sigma_k z, \quad z \sim \mathcal{N}(0, I) \quad (3)$$

115 where $\bar{\alpha}_k = \prod_{i=1}^k \alpha_i$ follows a predetermined noise schedule and σ_k controls the stochasticity of
 116 the reverse process. The forward (noising) process adds Gaussian noise, $x_k = \sqrt{\bar{\alpha}_k} a_0 + \sqrt{1 - \bar{\alpha}_k} \epsilon$,
 117 where a_0 is a clean action and $\epsilon \sim \mathcal{N}(0, I)$.

118 The key innovation of SDAC is a *Q-weighted* denoising loss that addresses a fundamental mismatch
 119 in applying diffusion models to RL standard denoising score matching fits the policy to whatever
 120 actions appear in the replay buffer, but in RL we want the policy to concentrate mass on *high-value*
 121 actions, not arbitrary ones. SDAC resolves this by reweighting the denoising loss by Q-values, so the
 122 generative capacity is steered toward high-reward modes of the action distribution. For each state s in
 123 the batch, SDAC (i) draws N_{mc} Monte Carlo action candidates $\{\hat{a}_i\}_{i=1}^{N_{mc}}$ from the current policy, (ii)
 124 scores each with the Q-function and converts the scores into softmax importance weights w_i , (iii)
 125 samples a diffusion step $k \sim \mathcal{U}\{1, \dots, T\}$ and i.i.d. noise $\epsilon^{(i)} \sim \mathcal{N}(0, I)$, and (iv) constructs the
 126 noisy targets $x_k^{(i)}$ via the forward diffusion process. Concretely,

$$w_i = \frac{\exp(c Q(s, \hat{a}_i) / \tau_{\text{ent}})}{\sum_{j=1}^M \exp(c Q(s, \hat{a}_j) / \tau_{\text{ent}})}, \quad x_k^{(i)} = \sqrt{\bar{\alpha}_k} \hat{a}_i + \sqrt{1 - \bar{\alpha}_k} \epsilon^{(i)}, \quad (4)$$

127

$$\mathcal{L}_{\text{SDAC}}(\theta) = \sum_{i=1}^{N_{\text{mc}}} w_i \left\| \epsilon_{\theta}(x_k^{(i)}, k, s) - \epsilon^{(i)} \right\|^2 \quad (5)$$

128 where $c = 5$ is the temperature scaling constant from the original work and $\bar{\alpha}_k$ is the cumulative noise
 129 schedule of Eq. 3. Each candidate \hat{a}_i thus serves as a clean-action target a_0 in the standard DDPM
 130 forward process, and the score-matching loss is evaluated on the candidate-specific noisy state $x_k^{(i)}$
 131 with its own injected noise $\epsilon^{(i)}$. The denominator in w_i is the entropy temperature itself; the original
 132 SDAC code expresses it as $\exp(\log \tau_{\text{ent}})$ because the log-temperature is the learned parameter, but
 133 the two forms are equivalent. At inference, N candidates are sampled and the one with the highest
 134 $\min(Q_{\phi_1}, Q_{\phi_2})$ is selected (best-of- N selection).

135 **Flow matching policies** Conditional flow matching [Lipman et al., 2023] offers an alternative genera-
 136 tive framework that replaces the complex noise schedule of diffusion with straight-line interpolation
 137 between noise and data. The forward process defines a linear path

$$x_t = (1 - t) a_0 + t \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \quad t \in [0, 1] \quad (6)$$

138 where $t = 0$ corresponds to the clean action and $t = 1$ to pure noise. A velocity network v_{θ} is trained
 139 to predict the conditional velocity field $v^* = \epsilon - a_0$ that transports points along these linear paths,

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t \sim \mathcal{U}[0,1], \epsilon \sim \mathcal{N}(0,I)} \left[\|v_{\theta}(x_t, t, s) - (\epsilon - a_0)\|^2 \right] \quad (7)$$

140 Sampling proceeds via Euler ODE integration from $t = 1$ (noise) to $t = 0$ (action), with the update
 141 $x_{t-\Delta t} = x_t - v_{\theta}(x_t, t, s) \cdot \Delta t$. Flow matching produces simpler velocity fields and more consistent
 142 candidates via deterministic ODE integration.

143 **Ensemble Exploration.** Bootstrapped DQN [Osband et al., 2016] maintains K Q-network heads
 144 $\{Q_{\phi_k}\}_{k=1}^K$, each trained on a bootstrapped data subset obtained by independently including each
 145 transition with probability p (Bernoulli masking). Thompson Sampling selects one head at the
 146 start of each episode and follows the greedy policy of throughout, producing *temporally coherent*
 147 exploration, the agent commits to one head’s strategy for the full episode rather than switching
 148 strategies every step, which enables deep exploration of long-horizon plans that single-step methods
 149 cannot. SUNRISE [Lee et al., 2021] replaces per-episode commitment with per-step UCB selection
 150 The shift to per-step selection trades temporal coherence for finer-grained uncertainty exploitation
 151 at every state, which works better in continuous-control settings where episodes are long and value
 152 uncertainty varies meaningfully across nearby states. REDQ [Chen et al., 2021] introduces random
 153 subset minimization for Bellman targets, with the target using $\min_{i \in \mathcal{M}} Q_{\phi_i}(s', a')$ where \mathcal{M} is a
 154 random subset of size M from the K Q-networks, providing a tunable pessimism-optimism knob
 155 ($M = 2$ approximates twin-Q, $M = K$ is maximally pessimistic).

156 3 Methodology

157 We saw in Section 1 that the naive extension of ensemble exploration to generative policies fails:
 158 K flow or diffusion heads trained against a shared Q-network collapse to functionally identical
 159 policies regardless of bootstrap masking or initialization, because the Q-weighted denoising gradient
 160 is head-invariant in expectation when the same Q drives every head’s importance weights. Prior work
 161 on ensemble exploration in deep RL [Osband et al., 2016, Lee et al., 2021, Chen et al., 2021] sidesteps
 162 this collapse only because Gaussian policies have a fundamentally different gradient structure: their
 163 parameters enter the loss directly through a reparametrized log-density, not indirectly through Q-
 164 weighted importance sampling. The transition to generative policies therefore requires a non-obvious
 165 architectural change.

166 ESAC makes this change. Each policy head is paired with its *own* independently initialized Q-network
 167 and trained with Q-weighted importance weights computed from that head’s Q rather than a shared
 168 one, so different Q-landscapes drive different policies even on shared data. A circular twin-Q pairing
 169 recovers SAC-style pessimism without doubling critic parameter count, and a trajectory-integrated
 170 disagreement bonus (TIDE) drives the resulting ensemble to act on its uncertainty during environment
 171 interaction. Figure 1 summarizes the full pipeline; the rest of this section describes the architecture
 172 (§3.1) and TIDE (§3.2) in detail.

173 3.1 Independent Actor-Critic Pairs

174 The core architectural decision in ESAC is to maintain K independently-initialized flow policy heads
 175 $\{\pi_{\theta_k}\}_{k=1}^K$ paired with K independently-initialized Q-networks $\{Q_{\phi_k}\}_{k=1}^K$. Twin-Q pessimism is

176 recovered through circular pairing where head k uses the index set $S_k = \{k, (k+1) \bmod K\}$ in its
 177 policy loss, so each Q_{ϕ_k} participates in exactly two pairings while keeping critic parameter count
 178 linear in K .

179 Each Q_{ϕ_k} is trained via Bellman backup with REDQ-style random subset minimization [Chen et al.,
 180 2021]. At every gradient step we redraw a random size- M subset $\mathcal{M}_t \subseteq \{1, \dots, K\}$ of the target
 181 Q-networks and use

$$y = r + \gamma \min_{i \in \mathcal{M}_t} Q_{\phi_i}(s', a'), \quad a' \sim \pi_{\theta_0}(\cdot | s') \quad (8)$$

182 with the same scalar target y training all K Q-networks at that step. Sampling a' from a single
 183 anchor head decouples policy diversity that is carried by the per-head importance weights below from
 184 value-target stability. Entropy regularization is realized implicitly in the actor through the Q-weighted
 185 softmax loss in Eq. 9, equivalent to maximizing $\mathbb{E}_{a \sim \pi}[Q] + \tau_{\text{ent}} \mathcal{H}[\pi]$ under the softmax distribution
 186 over the candidate pool. This avoids evaluating $\log \pi_{\theta}$ on the generative policy. Auto-tuning of τ_{ent}
 187 uses the standard SAC [Haarnoja et al., 2018] dual objective applied to the analytic Gaussian entropy
 188 of the action-time exploration noise $\mathcal{N}(0, (0.1 \tau_{\text{ent}})^2 I)$.

189 Per-head pairing is applied at the policy loss. Each head π_{θ_k} is trained with a Q-weighted denoising /
 190 flow-matching loss (Eq. 5, Eq. 7) using importance weights computed via its circular twin-Q pair,

$$w_{k,i} = \text{softmax}_i \left(\frac{c \cdot \min_{l \in S_k} Q_{\phi_l}(s, \hat{a}_i)}{\tau_{\text{ent}}} \right) \quad (9)$$

191 where the softmax is over the N_{mc} Monte Carlo candidates $\{\hat{a}_i\}_{i=1}^{N_{\text{mc}}}$ drawn from π_{θ_k} . For any $j \neq k$
 192 the index sets S_j and S_k overlap in at most one element, so each head’s gradient depends on at least
 193 one Q-network the other’s does not and the structural property formalized in Theorem 2. Since each
 194 Q_{ϕ_l} converges differently, the pairwise minima differ across heads, the weights $w_{k,i}$ differ across
 195 heads, and each policy is driven toward different regions of action space. The total policy loss is
 196 averaged across heads, $\mathcal{L}(\theta) = \frac{1}{K} \sum_{k=1}^K \mathcal{L}_k(\theta_k)$.

197 3.2 Trajectory Integrated Disagreement Exploration

198 UCB measures disagreement only at the final action, ignoring the iterative denoising structure.
 199 TIDE exploits this additional structure by measuring disagreement at multiple points along the
 200 denoising trajectory and using the integrated disagreement as the exploration bonus. Intuitively,
 201 when the K heads disagree throughout the generation process rather than only at the output, they
 202 represent genuinely distinct posterior samples and the associated state is a more informative one
 203 to visit. Concretely, for each candidate action a_j produced by head k and any denoising timestep
 204 $\tau \in \mathcal{T}_{\text{TIDE}} := \{0, T/4, T/2, 3T/4\}$, we evaluate all K velocity (or noise) networks on the noised
 205 intermediate $x_{\tau}^{(j)}$ and collect the prediction disagreement,

$$\Delta_{\tau}^{(j)} = \frac{1}{d_a} \sum_{d=1}^{d_a} \sqrt{\frac{1}{K-1} \sum_{k=1}^K \left(v_{\theta_k}(x_{\tau}^{(j)}, \tau, s)_d - \bar{v}_{\tau}(x_{\tau}^{(j)}, s)_d \right)^2} \quad (10)$$

206 where \bar{v}_{τ} is the mean velocity across heads (for diffusion, v_{θ_k} is replaced by the noise prediction ϵ_{θ_k}).
 207 The trajectory-integrated disagreement is

$$\widehat{\Delta}_{\text{traj}}^{(j)} = \sqrt{\frac{1}{|\mathcal{T}_{\text{TIDE}}|} \sum_{\tau \in \mathcal{T}_{\text{TIDE}}} \left(\Delta_{\tau}^{(j)} \right)^2}, \quad (11)$$

208 i.e. the root-mean-square disagreement over the sampled denoising timesteps for candidate j , used as
 209 the exploration bonus in action selection. This is the empirical estimator of the population quantity
 210 used in Proposition 1,

$$\Delta_{\text{traj}}^2 := \mathbb{E}_{t \sim \mathcal{T}} \left[\text{Var}_k f_{\theta_k} \left(x_t^{(k)}, t, s \right) \right], \quad (12)$$

211 where $x_t^{(k)}$ denotes the head- k -specific intermediate state at denoising time t . Since
 212 $(\Delta_{\tau}^{(j)})^2$ in Eq. (10) is already the per-timestep empirical variance across heads, $(\widehat{\Delta}_{\text{traj}}^{(j)})^2 =$
 213 $\frac{1}{|\mathcal{T}_{\text{TIDE}}|} \sum_{\tau} \text{Var}_k f_{\theta_k}(x_{\tau}^{(j)}, \tau, s)$ is an unbiased estimator of Δ_{traj}^2 when $x_{\tau}^{(j)}$ is drawn from the marginal
 214 of the reverse process. Candidates are then scored as

$$a_t^* = \arg \max_{a_j} \bar{Q}(s_t, a_j) + \beta(t) \cdot \widehat{\Delta}_{\text{traj}}^{(j)}, \quad (13)$$

215 where t denotes the environment step and a_t^* is the action executed at s_t . We use τ exclusively for the
 216 denoising timestep (Eq. 10) and t exclusively for the environment step (Eq. 13, Eq. 14) throughout
 217 this section.

218 For each environment step, scoring the M_c candidates per head requires evaluating all K velocity
 219 networks at each of the $|\mathcal{T}_{\text{TIDE}}| = 4$ timesteps on every noised intermediate, giving $|\mathcal{T}_{\text{TIDE}}| \cdot K \cdot$
 220 $(K \cdot M) = 4K^2M$ velocity forward passes for exploration scoring alone. Empirically, because each
 221 velocity network is a small MLP and the $4K^2M$ evaluations are batched into a single forward call,
 222 TIDE inference adds $\approx 35\%$ to per step wall clock time over best-of- N Q-selection at training time
 223 and is disabled at evaluation.

224 A key failure mode of any disagreement-based bonus is that in relatively unimodal environments
 225 (e.g. HalfCheetah), the K velocity networks may disagree simply due to parameter noise, even
 226 when the Q-landscape is flat. We provide empirical evidence that Ant and Humanoid possess
 227 relatively large multimodal Q-landscapes supporting multiple qualitatively distinct high-reward action
 228 strategies, while HalfCheetah exhibits a comparatively unimodal structure, via Q-value distribution
 229 analysis and behavioural clustering of independently trained agents across all three environments; We
 230 provide empirical evidence in Appendix E that Ant and Humanoid possess multimodal Q-landscapes
 231 supporting multiple qualitatively distinct high-reward action strategies, while HalfCheetah exhibits a
 232 comparatively unimodal structure.

$$\beta(t) = \beta_0 \cdot \max(0, \rho_t), \quad \rho_t = \frac{\text{Cov}(\Delta_{\text{traj}}, \sigma_Q)}{\sigma(\Delta_{\text{traj}}) \sigma(\sigma_Q)}. \quad (14)$$

233 When velocity disagreement and Q-uncertainty agree on where the agent is uncertain ($\rho_t \approx 1$), the
 234 full bonus is applied, and when they disagree ($\rho_t \leq 0$) the bonus is suppressed and TIDE falls back
 235 to pure exploitation at that step.

236 4 Theoretical Analysis

237 We provide formal justification for three claims.

238 **Theorem 1** (Collapse under Shared Backbone and Shared Twin-Q). *Let K policy heads share a*
 239 *backbone ϕ and a single twin-Q pair, with per-head Bernoulli masks $m_{k,i} \sim \text{Bernoulli}(p)$. Then the*
 240 *expected gradient is head-invariant:*

$$\mathbb{E}_{m_k}[\nabla_{\eta_k} \mathcal{L}_k] = p \cdot \nabla_{\eta_k} \mathcal{L}^{(0)}(\theta_k; Q_\phi) \quad \forall k, \quad (15)$$

241 *and the pairwise function-space distance $\|f_{\theta_i}(\cdot) - f_{\theta_j}(\cdot)\|^2$ contracts to zero at rate $O(\eta_t)$ under*
 242 *standard SGD. Bootstrap masking cannot prevent collapse since $\text{Var}[\bar{m}_k] = p(1-p)/B$ is a factor*
 243 *$1/B$ smaller than the per-sample variant of Bootstrapped DQN [Osband et al., 2016].*

244 Figure 4 in Appendix D confirms collapse empirically the off-diagonal action-cosine mean reaches
 245 exactly 1.000 by 200k steps under shared Q , while under independent Q it stabilizes near 0.68.

246 **Lemma 1** (Non-identifiability of Heads under Shared Q). *$\mathcal{L}^{(0)}$ is invariant to permutations of the*
 247 *head index, so the K -head model is non-identifiable up to head reparameterization. Function-space*
 248 *collapse therefore does not require parameter-space collapse; heads can retain distinct parameters*
 249 *that all implement the same function.*

250 Figure 2 confirms this that under shared Q function-space distance crashes to near zero by 200k steps
 251 while parameter-space distance remains substantially nonzero throughout. More detailed explanation
 252 on this is in Appendix D

253 **Theorem 2** (Initial-Step Gradient Non-Collapse under Circular Twin-Q Pairing). *Replace the shared*
 254 *twin-Q with K independently-initialized Q-networks $Q_{\phi_k} \sim \mathcal{N}(0, \sigma_{\text{init}}^2 I)$ and use circular twin-Q*
 255 *pairing $S_k = \{k, (k+1) \bmod K\}$ for per-head importance weights. For any $j \neq k$ the index sets S_k*
 256 *and S_j overlap in at most one element, so each head’s gradient depends on at least one Q-network*
 257 *the other’s does not. Under Assumptions 1–2,*

$$\mathbb{E} \|\nabla_{\eta_k} \mathcal{L}_k - \nabla_{\eta_j} \mathcal{L}_j\|^2 \geq c \cdot \frac{\sigma_{\text{init}}^2}{K} \quad \text{for all } j \neq k, \quad (16)$$

258 *breaking the head-invariance that drives collapse in Theorem 1.*

259 **Proposition 1** (TIDE as a Tighter Epistemic Uncertainty Proxy). *Let $\Delta_0 := \text{std}_k f_{\theta_k}(x_0, 0, s)$ be ter-*
 260 *minial disagreement and $\Delta_{\text{traj}}^2 := \mathbb{E}_{t \sim \mathcal{T}} [\text{Var}_k f_{\theta_k}(x_t^{(k)}, t, s)]$ be trajectory-integrated disagreement*
 261 *(Eq. 12). Under Assumption 1,*

$$\text{Var}_k[a_k(s)] \leq C_{L_x}^2 \cdot \left(\sum_{t \in \mathcal{T}} L_g(t) \right)^2 \cdot \Delta_{\text{traj}}^2 + O(\Delta_0^2), \quad (17)$$

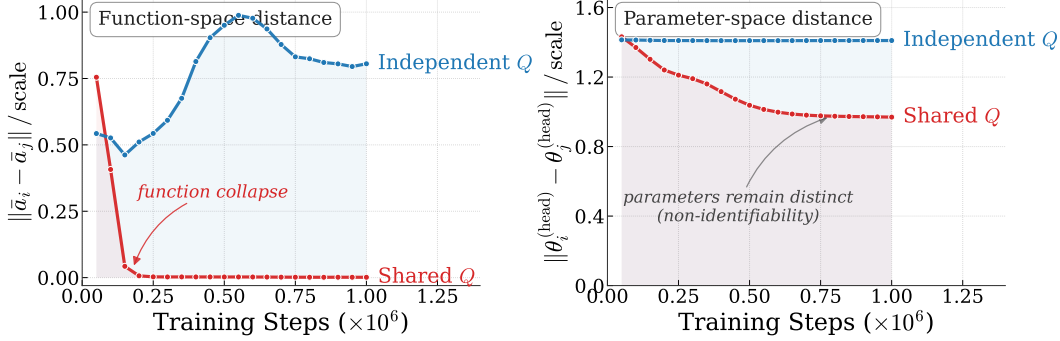


Figure 2: Empirical confirmation of Lemma 1 on Ant-v4 with $K=5$. Each panel reports a mean pairwise inter-head L2 distance, normalized by the mean L2 magnitude of a single head’s output. Curves show shared Q (red) and independent Q (blue); annotations highlight the collapse and non-identifiability regimes.

262 *so inter-head action variance is bounded by trajectory-integrated disagreement alone, while Δ_0*
 263 *captures only the final denoising step. On Ant-v4, Δ_{traj} is up to $3\times$ larger than Δ_0 during mid-*
 264 *training, providing a richer signal exactly when the policy has begun to specialize but not yet*
 265 *committed to a single strategy.*

266 Full proofs, figures, and the Observation on finite-horizon diversity preservation are in Appendix D.

267 5 Results and Discussions

268 We evaluate ESAC on standard MuJoCo continuous control benchmarks, comparing against the
 269 SDAC baseline [Ma et al., 2025] and other baselines. Our central question is whether directed,
 270 posterior-driven exploration meaningfully outperforms the noise-based exploration used by every
 271 prior generative-policy RL method on tasks where exploration is the actual bottleneck. We therefore
 272 expect ESAC’s gains to be largest on tasks like Ant-v4 and Humanoid-v4 with multimodal action
 273 distributions and tight termination conditions, and smallest on smooth landscapes like HalfCheetah
 274 where any reasonable exploration strategy suffices.

275 5.1 Setup

276 We evaluate on Mujoco environments HalfCheetah, Hopper, Walker2d, Ant, Pusher, Reacher, and
 277 Humanoid running all experiments for 1M environment steps except Humanoid which runs for 5M,
 278 with data collected through 5 vectorized environments. All experiments and training is carried out
 279 on NVIDIA RTX A6000 48GB VRAM and methods share the same hyperparameters, including a
 280 learning rate of 3×10^{-4} with linear decay to 3×10^{-5} , an α learning rate of 7×10^{-3} , $T = 20$
 281 denoising or flow steps, batch size 256, a replay buffer of 10^6 transitions, $\tau = 0.005$, $\gamma = 0.99$, loss
 282 scaling probability $p = 0.8$, and REDQ $M = 2$. Full details are provided in Appendix F.

283 5.2 Results

284 Figure 3 shows the training learning curves across the eight MuJoCo locomotion tasks of our
 285 benchmark. We compare ESAC with the flow-matching and diffusion backbones (red and green)
 286 against the SDAC baseline (orange), DPMD (blue), and the SAC baseline (gray). ESAC with the
 287 flow-matching backbone matches or exceeds the strongest baseline on Ant-v4, Walker2d-v4, and
 288 Humanoid-v4, the three tasks where directed exploration matters most, and stays competitive on
 289 the comparatively unimodal tasks Hopper-v4 and HalfCheetah-v4. Table 1 reports final evaluation
 290 returns across all environments. ESAC (FM) achieves the best result on three of the five locomotion
 291 tasks shown in the figure, with Ant-v4 (5646 ± 1001), Walker2d-v4 (4305 ± 529), and Humanoid-
 292 v4 (6773 ± 510), and remains within $\sim 15\%$ of the leader on the remaining two. Relative to the
 293 SDAC backbone [Ma et al., 2025], ESAC (FM) yields a $7.08\times$ improvement on Ant-v4, $1.36\times$ on
 294 Walker2d-v4, $1.18\times$ on Humanoid-v4, and $1.12\times$ on Hopper-v4.

295 The improvement scales monotonically with the exploration difficulty of the task. Ensemble diversity
 296 helps most when the Q -landscape is multimodal and high-dimensional, and yields diminishing returns
 297 when a single mode suffices. The signature is visible in the saturation pattern on HalfCheetah-v4,
 298 where the four SDAC-family methods all converge to within $\sim 14\%$ of each other (10,328–11,949),
 299 while the gap between the best and worst diffusion-policy method on Ant exceeds $58\times$ (95 for DPPO
 300 vs 5646 for ESAC FM). Against the strongest single competing baseline per environment, ESAC
 301 (FM) outperforms DPMD [Ma et al., 2025] on Ant (5646 vs 5520), Walker2d (4305 vs 4280), and

302 Humanoid (6773 vs 6710); DPMD edges out on Hopper-v4 (3220 vs 2759) and the original SDAC
 303 retains a slim HalfCheetah-v4 lead. On Humanoid-v4, ESAC (FM) also matches the variance profile
 304 of the strongest stable diffusion-policy baselines its seed standard deviation (± 510) is an order of
 305 magnitude tighter than the deeper diffusion baselines that exhibit seed-level training collapses on 17-
 306 DoF control (DACER ± 3210 , SAC ± 2415) and the parameter matched comparison in Appendix H
 307 shows SDAC and SAC scaled to ESAC’s budget still trail ESAC (FM) by roughly 1000 return points
 308 on Ant-v4, confirming that the gains stem from ensemble structure rather than added capacity or
 309 candidate count.

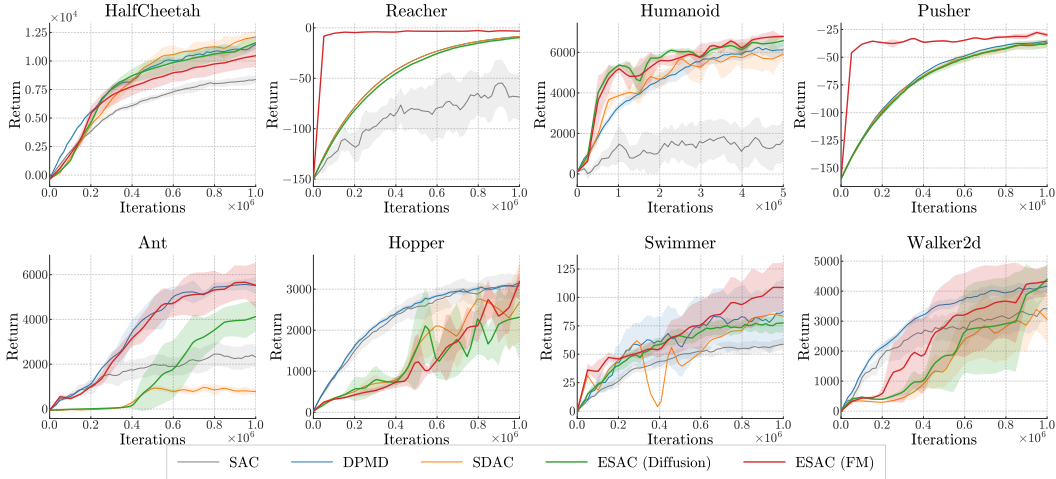


Figure 3: Training learning curves on five MuJoCo locomotion tasks trained for 1M steps (200k iterations across 5 vectorized environments) for top 5 baselines

Table 1: Final evaluation return on MuJoCo locomotion tasks, mean \pm std over 5 seeds. **Bold** marks methods within one standard deviation of the highest mean per environment.

Method		HALFCHEETAH	REACHER	HUMANOID	PUSHER
Classic Model-Free RL	PPO [Schulman et al., 2017]	4715 \pm 810	-7.56 \pm 10.68	1015 \pm 280	-27.52 \pm 3.40
	TD3 [Fujimoto et al., 2018]	8345 \pm 645	-3.19 \pm 0.04	5705 \pm 405	-25.23 \pm 0.96
	SAC [Haarnoja et al., 2018]	8745 \pm 415	-61.58 \pm 58.11	2945 \pm 2415	-34.28 \pm 1.16
Diffusion Policy RL	QSM [Psenka et al., 2024]	10495 \pm 510	-4.45 \pm 0.45	5520 \pm 480	-83.18 \pm 2.20
	DIPO [Yang et al., 2023]	9285 \pm 720	-3.28 \pm 0.03	4790 \pm 985	-32.11 \pm 0.44
	DACER [Wang et al., 2024]	11055 \pm 285	-3.12 \pm 0.07	2880 \pm 3210	-31.02 \pm 0.18
	QVPO [Ding et al., 2024]	7195 \pm 1015	-35.89 \pm 13.25	445 \pm 95	-134.21 \pm 1.08
	DPPO [Ren et al., 2024]	1245 \pm 415	-5.62 \pm 1.67	510 \pm 75	-95.45 \pm 11.87
	DPMD [Ma et al., 2025]	11750 \pm 670	-3.18 \pm 0.08	6710 \pm 525	-30.18 \pm 0.33
SDAC [Ma et al., 2025]	11949 \pm 659	-3.43 \pm 0.41	5766 \pm 386	-32.28 \pm 4.98	
Ours	ESAC (Diff)	11310 \pm 728	-3.84 \pm 1.96	6493 \pm 672	-31.52 \pm 5.88
	ESAC (FM)	10328 \pm 1119	-3.03 \pm 0.87	6773 \pm 510	-28.90 \pm 4.53
Method		ANT	HOPPER	SWIMMER	WALKER2D
Classic Model-Free RL	PPO [Schulman et al., 2017]	3320 \pm 945	3155 \pm 220	84.5 \pm 12.4	4220 \pm 740
	TD3 [Fujimoto et al., 2018]	3690 \pm 1180	1855 \pm 985	71.9 \pm 15.3	2640 \pm 1190
	SAC [Haarnoja et al., 2018]	2480 \pm 820	3110 \pm 345	63.5 \pm 10.2	3110 \pm 950
Diffusion Policy RL	QSM [Psenka et al., 2024]	870 \pm 195	2715 \pm 510	57.0 \pm 7.7	2480 \pm 805
	DIPO [Yang et al., 2023]	1020 \pm 145	1245 \pm 690	46.7 \pm 2.9	2050 \pm 1380
	DACER [Wang et al., 2024]	4080 \pm 615	3145 \pm 145	103.0 \pm 45.8	3275 \pm 1640
	QVPO [Ding et al., 2024]	745 \pm 290	2810 \pm 555	53.4 \pm 5.0	2415 \pm 1095
	DPPO [Ren et al., 2024]	95 \pm 35	2245 \pm 580	106.1 \pm 6.5	1085 \pm 720
	DPMD [Ma et al., 2025]	5520 \pm 245	3220 \pm 85	84.33 \pm 34.56	4280 \pm 320
SDAC [Ma et al., 2025]	798 \pm 190	2472 \pm 993	83.84 \pm 16.9	3176 \pm 782	
Ours	ESAC (Diff)	4043 \pm 617	2284 \pm 1147	77.05 \pm 16.04	4172 \pm 483
	ESAC (FM)	5646 \pm 1001	2759 \pm 437	107.83 \pm 23.82	4305 \pm 529

310 6 Related Work

311 **Generative policies for RL.** Diffusion Policy [Chi et al., 2023] pioneered diffusion models for
312 behavioral cloning in robotics. For online RL, SDAC [Ma et al., 2025] introduced Q-weighted
313 score matching for diffusion actors, QVPO [Ding et al., 2024] combined Q-weighted diffusion with
314 variational policy optimization, and Q-score matching [Psenka et al., 2024] connected the policy score
315 to the Q-function gradient. FPO [McAllister et al., 2025] introduced flow matching [Lipman et al.,
316 2023, Liu et al., 2023] for policy gradients. In the offline setting, diffusion policies have been applied
317 via implicit Q-learning [Hansen-Estruch et al., 2024], conditional generation [Wang et al., 2023], and
318 out-of-distribution generalization [Ada et al., 2024], while policy-guided diffusion [Jackson et al.,
319 2024] fine-tunes pretrained diffusion policies via RL objectives. All of these methods, both online
320 and offline, use a single generative policy with noise-based exploration. ESAC is the first to extend
321 ensemble exploration to generative policy classes, and the first to identify and resolve the shared-Q
322 collapse mode that breaks the naive ensemble extension.

323 **Ensemble methods in RL.** The use of ensembles for exploration traces back to Thompson Sam-
324 pling [Thompson, 1933, Russo et al., 2018]. Bootstrapped DQN [Osband et al., 2016] and randomized
325 prior functions [Osband et al., 2018] introduced multi-head Q-networks with Thompson Sampling for
326 deep exploration in discrete domains. REDQ [Chen et al., 2021] scaled Q-ensembles to continuous
327 control with high update-to-data ratios. SUNRISE [Lee et al., 2021] combined ensemble disagreement
328 with UCB-style optimism [Auer et al., 2002]. OAC [Ciosek et al., 2019] used approximate upper
329 confidence bounds for Gaussian policies. EDAC [An et al., 2021] penalized Q-function uncertainty
330 for offline RL, and Dropout Q-functions [Hiraoka et al., 2022] provided a lightweight alternative
331 to full ensembles. TEEN [Li et al., 2023] explicitly promoted trajectory diversity across ensemble
332 members. Theoretical foundations include tighter Bayesian regret bounds [Moradipari et al., 2023]
333 and approximate posterior sampling [Ishfaq et al., 2024]. All of these methods operate on Gaussian or
334 otherwise unimodal policy classes, and our results (Theorem 1) show that their default shared-Q de-
335 signs do not transfer to generative policies without the architectural change ESAC introduces. We also
336 note a surface resemblance between ESAC’s ensemble-mean aggregation and value-decomposition
337 methods such as VDN [Sunehag et al., 2018] and QMIX [Rashid et al., 2020], where the total Q
338 value factorizes a joint Q-function across cooperative agents with private observations. The two
339 are structurally distinct, ESAC’s K Q-networks are independent estimates of the *same* single-agent
340 Q-function on the *same* observation, summarized by an ensemble mean, whereas VDN’s Q_i are
341 pieces of a joint Q-function across *different* agents. The shared shape (a linear combination of K
342 Q-values) addresses orthogonal problems.

343 **Exploration in continuous control.** Beyond ensembles, exploration has been addressed via curiosity-
344 driven intrinsic motivation [Pathak et al., 2017], random network distillation [Burda et al., 2019],
345 count-based bonuses [Bellemare et al., 2016], and maximum entropy objectives [Harnoja et al.,
346 2018]. These methods target a different exploration regime they augment the reward function with a
347 state-novelty signal that depends on visitation history, encoding-prediction error, or policy entropy
348 none of which use the structure of an ensemble’s posterior disagreement to direct exploration. They
349 are complementary to ESAC rather than competing alternatives, and could in principle be added on
350 top of our ensemble framework as additional reward shaping. We do not benchmark against them here
351 because they address state-level novelty (where in S the agent has not been) rather than action-level
352 posterior diversity (which actions in A the agent is uncertain about), which is the regime ESAC is
353 designed for. We treat their integration as future work.

354 7 Conclusion

355 We explored whether ensemble exploration, a proven technique for simple Gaussian policies, transfers
356 to the increasingly popular family of generative RL policies. Our central finding is that it does not, at
357 least not without independent Q-targets. The Q-weighted denoising loss that trains generative policies
358 creates a gradient alignment failure mode where shared Q-functions drive all ensemble heads to
359 identical policies, regardless of bootstrap masking or architectural independence. With independent
360 actor-critic pairs and our new TIDE exploration rule, ESAC delivers substantial gains that scale with
361 task difficulty, especially on tasks that demand exploration the most. The gains are regime-dependent:
362 Appendix I shows ESAC plateaus on unimodal manipulation skills where coverage helps less than
363 sharpness, pointing to multi-solution manipulation and adaptive K as natural future work, alongside
364 extending ESAC to offline pretraining and offline-to-online finetuning where Q-uncertainty behaves
365 differently. One limitation is that ESAC trades a $3.4\times$ parameter and $\sim 35\%$ wall-clock overhead for
366 these gains.

367 References

- 368 Suzan Ece Ada, Erhan Oztop, and Emre Ugur. Diffusion policies for out-of-distribution generalization in offline
369 reinforcement learning. *IEEE Robotics and Automation Letters*, 9(4):3116–3123, 2024.
- 370 Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline reinforcement
371 learning with diversified Q-ensemble. In *Advances in Neural Information Processing Systems*, 2021.
- 372 Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem.
373 *Machine Learning*, 47(2-3):235–256, 2002.
- 374 Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying
375 count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems*,
376 2016.
- 377 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George
378 Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable
379 transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- 380 Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. In
381 *International Conference on Learning Representations*, 2019.
- 382 Xinyue Chen, Che Wang, Zijian Zhou, and Keith Ross. Randomized ensembled double Q-learning: Learning
383 fast without a model. In *International Conference on Learning Representations*, 2021.
- 384 Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song.
385 Diffusion policy: Visuomotor policy learning via action diffusion. In *Robotics: Science and Systems*, 2023.
- 386 Kamil Ciosek, Quan Vuong, Robert Loftin, and Katja Hofmann. Better exploration with optimistic actor-critic.
387 In *Advances in Neural Information Processing Systems*, 2019.
- 388 Shutong Ding, Ke Zheng, Chi Zhang, Jia Zhao, Yiping Shi, and Chao Chen. Diffusion-based reinforcement
389 learning via Q-weighted variational policy optimization. In *Advances in Neural Information Processing*
390 *Systems*, 2024.
- 391 Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic
392 methods. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, 2018.
- 393 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum
394 entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*,
395 2018.
- 396 Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine. IDQL:
397 Implicit Q-learning as an actor-critic method with diffusion policies. In *International Conference on Machine*
398 *Learning*, 2024.
- 399 Takuya Hiraoka, Takahisa Imagawa, Taisei Hashimoto, Takahiro Onishi, and Yoshimasa Tsuruoka. Dropout
400 Q-functions for doubly efficient reinforcement learning. In *International Conference on Learning Representa-*
401 *tions*, 2022.
- 402 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural*
403 *Information Processing Systems*, 2020.
- 404 Haque Ishfaq, Yixin Tan, Yu Yang, Qingfeng Lan, Jianfeng Lu, A Rupam Mahmood, Doina Precup, and Pan Xu.
405 More efficient randomized exploration for reinforcement learning via approximate sampling. *Reinforcement*
406 *Learning Journal*, 1:1558–1595, 2024. URL [https://rlj.cs.umass.edu/2024/papers/Paper148.](https://rlj.cs.umass.edu/2024/papers/Paper148.html)
407 [html](https://rlj.cs.umass.edu/2024/papers/Paper148.html).
- 408 Matthew T Jackson, Michael T Matthews, Chris Lu, Benjamin Ellis, Shimon Whiteson, and Jakob Foerster.
409 Policy-guided diffusion. *Reinforcement Learning Journal*, 2024.
- 410 Kimin Lee, Michael Laskin, Aravind Srinivas, and Pieter Abbeel. SUNRISE: A simple unified framework for
411 ensemble learning in deep reinforcement learning. In *International Conference on Machine Learning*, 2021.
- 412 Chao Li, Chen Gong, Qiang He, and Xinwen Hou. Keep various trajectories: Promoting exploration of ensemble
413 policies in continuous control. In *Advances in Neural Information Processing Systems*, volume 36, pages
414 5223–5235, 2023.
- 415 Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver,
416 and Daan Wierstra. Continuous control with deep reinforcement learning. *International Conference on*
417 *Learning Representations*, 2016.
- 418 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative
419 modeling. In *International Conference on Learning Representations*, 2023.
- 420 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data
421 with rectified flow. In *International Conference on Learning Representations*, 2023.

- 422 Haitong Ma, Hengkai Yan, Jianye Hao, and Zhen Xu. Efficient online reinforcement learning for diffusion
423 policy. In *International Conference on Machine Learning*, 2025.
- 424 David McAllister, Songwei Ge, Brent Yi, Chung Min Kim, Ethan Weber, Hongsuk Choi, Haiwen Feng, and
425 Angjoo Kanazawa. Flow matching policy gradients. *arXiv preprint arXiv:2507.21053*, 2025.
- 426 Ahmadreza Moradipari, Mohammad Pedramfar, Thodoris Lykouris, and Richard Combes. Improved Bayesian
427 regret bounds for Thompson sampling in reinforcement learning. In *Advances in Neural Information*
428 *Processing Systems*, 2023.
- 429 Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped
430 DQN. In *Advances in Neural Information Processing Systems*, 2016.
- 431 Ian Osband, John Aslanides, and Albin Cassirer. Randomized prior functions for deep reinforcement learning.
432 In *Advances in Neural Information Processing Systems*, 2018.
- 433 Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-
434 supervised prediction. In *International Conference on Machine Learning*, 2017.
- 435 Michael Psenka, Alejandro Escontrela, Pieter Abbeel, and Yi Ma. Learning a diffusion model policy from
436 rewards via Q-score matching. In *International Conference on Machine Learning*, 2024.
- 437 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder de Witt, Gregory Farquhar, Jakob Foerster, and Shimon
438 Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of*
439 *Machine Learning Research*, 21(178):1–51, 2020.
- 440 Allen Z. Ren, Justin Lidard, Lars L. Ankile, Anthony Simeonov, Pulkit Agrawal, Anirudha Majumdar, Ben-
441 jamin Burchfiel, Hongkai Dai, and Max Simchowitz. Diffusion policy optimization. *arXiv preprint*
442 *arXiv:2409.00588*, 2024.
- 443 Daniel J Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen. A tutorial on Thompson
444 sampling. *Foundations and Trends in Machine Learning*, 11(1):1–96, 2018.
- 445 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization
446 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 447 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole.
448 Score-based generative modeling through stochastic differential equations. *International Conference on*
449 *Learning Representations*, 2021.
- 450 Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg,
451 Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-decomposition
452 networks for cooperative multi-agent learning based on team reward. In *Proceedings of the 17th International*
453 *Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*, pages 2085–2087, 2018.
- 454 William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence
455 of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- 456 Yinuo Wang, Likun Wang, Yuxuan Jiang, Wenjun Zou, Tong Liu, Xujie Song, Wenxuan Wang, Liming Xiao,
457 Jiang Wu, Jingliang Duan, and Shengbo Eben Li. Diffusion actor-critic with entropy regulator. In *Advances*
458 *in Neural Information Processing Systems (NeurIPS)*, 2024.
- 459 Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy class for
460 offline reinforcement learning. In *International Conference on Learning Representations*, 2023.
- 461 Long Yang, Zhixiong Huang, Fenghao Lei, Yucun Zhong, Yiming Yang, Cong Fang, Shiting Wen, Binbin Zhou,
462 and Zhouchen Lin. Policy representation via diffusion probability model for reinforcement learning. *arXiv*
463 *preprint arXiv:2305.13122*, 2023.
- 464 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine.
465 Meta-World: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on*
466 *Robot Learning*, 2020.

467 **A Proof of Theorem 1 (Collapse under Shared Backbone and Shared Twin-Q)**

468 *Proof.* Recall the deployed Q-weighted loss for head k under the shared-backbone, shared-twin-Q
469 regime,

$$\mathcal{L}_k(\theta_k; Q_\phi) = \bar{m}_k \cdot \sum_{i=1}^B w_i(Q_\phi) \cdot \ell_i(\theta_k), \quad (18)$$

470 where $\bar{m}_k = \frac{1}{B} \sum_i m_{k,i}$ is the per-head Bernoulli-mask mean and $w_i(Q_\phi) =$
471 $\text{softmax}(\min(Q_{\phi_1}, Q_{\phi_2})(s_i, \hat{a}_i) \cdot c / \exp(\alpha))$ depends only on the shared twin-Q pair, hence is
472 head-invariant. The per-sample loss is $\ell_i(\theta_k) = \|f_{\theta_k}(x_t^{(i)}, t, s_i) - \text{target}_i\|^2$ with $f_{\theta_k} = h_k \circ \phi$
473 (shared backbone ϕ , head-specific linear h_k).

474 **Part (i): Expected gradient direction is head-invariant.** The mask is independent of the data and
475 parameters, so

$$\mathbb{E}_{m_k}[\nabla_{\eta_k} \mathcal{L}_k] = \mathbb{E}[\bar{m}_k] \cdot \sum_i w_i(Q_\phi) \cdot \nabla_{\eta_k} \ell_i(\theta_k) \quad (19)$$

$$= p \cdot \nabla_{\eta_k} \mathcal{L}^{(0)}(\theta_k; Q_\phi), \quad (20)$$

476 where $\mathcal{L}^{(0)} = \sum_i w_i \ell_i$ is the unmasked Q-weighted loss. The right-hand side depends on θ_k but
477 not on the head index k itself, and for any two heads i, j at the same parameter point the expected
478 gradient is identical. The same reasoning applies to the shared-backbone gradient $\nabla_\xi \mathcal{L}_k$, where
479 contributions from all K heads sum into the same trunk parameters with head-invariant magnitude.

480 **Part (ii): Function-space collapse.** Assume ℓ_i is L -smooth as a function of the network output
481 $f_{\theta_k}(x, t, s)$. Define the per-head output function $F_k(\cdot) := f_{\theta_k}(\cdot)$ and the function-space discrepancy
482 $d_t := \|F_i^{(t)}(\cdot) - F_j^{(t)}(\cdot)\|_\mu^2$, where $\|\cdot\|_\mu$ denotes the L^2 norm over the candidate-action distribution
483 μ . Because the expected gradient $\mathbb{E}_{m_k}[\nabla_{\eta_k} \mathcal{L}_k] = p \nabla_{\eta_k} \mathcal{L}^{(0)}(\theta_k; Q_\phi)$ is identical for all heads k
484 (from Part (i)), and $\mathcal{L}^{(0)}$ depends on θ_k only through the output F_k , the SGD iterates for F_i and F_j
485 satisfy the same expected update rule. By L -smoothness applied in output space and the Polyak–
486 Łojasiewicz (PL) condition on $\mathcal{L}^{(0)}$ viewed as a functional of F_k (which holds without requiring a
487 unique minimiser in parameter space, in contrast to strong convexity), we obtain

$$d_{t+1} \leq (1 - 2\eta_t p \mu_{\text{PL}}) d_t + O(\eta_t^2), \quad (21)$$

488 where $\mu_{\text{PL}} > 0$ is the PL constant of $\mathcal{L}^{(0)}$ in function space. Under standard step-size schedules
489 $\sum_t \eta_t = \infty$, $\sum_t \eta_t^2 < \infty$, this gives $d_t \rightarrow 0$ at rate $O(\eta_t)$, i.e. $\|F_i^{(t)}(\cdot) - F_j^{(t)}(\cdot)\|_\mu^2 \rightarrow 0$.

490 *Note:* $\mathcal{L}^{(0)}$ is permutation-invariant (Lemma 1), so its minimisers form a manifold rather than a
491 unique point; the parameter-space distance $\|\eta_i - \eta_j\|$ need not vanish and empirically does not.
492 The collapse is purely in function space: all heads converge to the same input–output map while
493 potentially retaining distinct parameter vectors on the permutation-symmetry orbit.

494 **Part (iii): Variance bound.** The Bernoulli-mask variance of the deployed scalar mask is $\text{Var}[\bar{m}_k] =$
495 $\text{Var}[\frac{1}{B} \sum_i m_{k,i}] = p(1-p)/B$. The per-sample mask used by Bootstrapped DQN [Osband et al.,
496 2016] contributes $p(1-p) \sum_i w_i^2 (\nabla \ell_i)^2$, which is a factor of B larger but still vanishes in expectation.
497 Either form is dominated by the head-invariant mean signal of part (i) and cannot break the collapse.

498 **Remark.** Even with randomized prior functions [Osband et al., 2018] added to each head’s output,
499 the gradient of the *trainable* parameters still follows the shared Q-landscape, so priors alone cannot
500 resolve the collapse. \square

501 **B Proof of Theorem 2 (Initial-Step Gradient Non-Collapse under Circular**
502 **Twin-Q Pairing)**

503 **Assumption 1** (Regularity). *The MDP $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$ has bounded rewards $|r(s, a)| \leq R_{\text{max}}$ and*
504 *the Q-function class $\mathcal{F} = \{Q_\phi : \phi \in \Phi\}$ has bounded complexity (e.g., bounded Rademacher*
505 *complexity).*

506 **Assumption 2** (Approximation Diversity). *There exist states (s, a) with positive function approx-*
507 *imation uncertainty $\sigma_{\text{approx}}^2(s, a) := \text{Var}_{\mathcal{D} \sim P}[\hat{Q}(s, a; \mathcal{D})] > 0$, where $\hat{Q}(s, a; \mathcal{D})$ is the Q-function*
508 *estimate trained on dataset \mathcal{D} .*

509 *Proof.* The argument is purely structural and holds at $t = 0$ under independent initialization. Write
 510 $S_k = \{k, (k+1) \bmod K\}$ for head k 's pessimism index set, so the importance weights take the
 511 form $w_{k,i} \propto \exp(\beta q_{k,i})$ with $q_{k,i} := \min_{\ell \in S_k} Q_{\phi_\ell}(s, \hat{a}_i)$, and the per-head policy gradient is
 512 $\nabla_{\eta_k} \mathcal{L}_k = \sum_{i=1}^M w_{k,i} \nabla_{\eta_k} \ell_i(\eta_k)$ where ℓ_i is the per-sample denoising or velocity loss on candidate
 513 \hat{a}_i .

514 The geometry of the circular pairing is what drives the result. For any two distinct heads $j \neq k$
 515 (with $K \geq 3$), the index sets S_k and S_j each consist of two consecutive integers modulo K , so
 516 they coincide only when $j \equiv k \pmod{K}$ and otherwise overlap in at most one element, as can be
 517 checked case by case for $j = k \pm 1, k \pm 2$. Consequently, at least one critic index $\ell^* \in S_k \setminus S_j$
 518 always exists, and head k 's importance weights depend on ϕ_{ℓ^*} while head j 's do not.

519 This single non-overlapping critic is sufficient to break gradient symmetry. Holding all ϕ_ℓ with
 520 $\ell \neq \ell^*$ fixed and viewing $q_{k,i}$ as a function of ϕ_{ℓ^*} alone, independent Gaussian initialization
 521 $\phi_{\ell^*} \sim \mathcal{N}(0, \sigma_{\text{init}}^2 I)$ makes $Q_{\phi_{\ell^*}}(s, \hat{a}_i)$ a random variable with positive variance $\Theta(\sigma_{\text{init}}^2)$ (the implicit
 522 constant absorbing the layerwise gain). With probability strictly greater than $1/2$, this random
 523 Q-value is the active argument of the min over S_k , so $\partial q_{k,i} / \partial \phi_{\ell^*}$ is non-zero on a set of positive
 524 measure. By contrast, $q_{j,i}$ has no dependence on ϕ_{ℓ^*} at all because $\ell^* \notin S_j$, so head j 's weights are
 525 deterministically constant in this direction.

526 It remains to translate this asymmetry into the importance weights into a lower bound on the gradient
 527 gap. Let $u_i := \nabla_{\eta_k} \ell_i(\eta_k)$ and $v_i := \nabla_{\eta_j} \ell_i(\eta_j)$ denote the head- k and head- j gradients of the
 528 same per-sample loss evaluated at their respective initial parameters. Assumption 2 guarantees
 529 $\|u_i\|, \|v_i\| \geq c_\ell > 0$ in expectation over the candidate distribution, and the independent initialization
 530 of η_k and η_j implies $\mathbb{E}[\langle u_i, v_i \rangle] = 0$ since the directions $u_i / \|u_i\|$ and $v_i / \|v_i\|$ are independent on the
 531 unit sphere. Expanding the squared gradient difference and dropping the cross term yields

$$\mathbb{E} \left\| \nabla_{\eta_k} \mathcal{L}_k - \nabla_{\eta_j} \mathcal{L}_j \right\|^2 \geq \mathbb{E} \left[\sum_i (w_{k,i} - w_{j,i})^2 \|u_i\|^2 \right] + \mathbb{E} \left[\sum_i w_{j,i}^2 \|u_i - v_i\|^2 \right]. \quad (22)$$

532 The first term inherits a lower bound from the asymmetry above. Along the direction ϕ_{ℓ^*} , the weights
 533 $w_{k,i}$ vary through the softmax composed with min while $w_{j,i}$ does not, and the candidate Q-values
 534 are bounded by Assumption 1, so the softmax does not saturate. Hence there exists $c_w > 0$ depending
 535 only on the temperature β and on R_{max} such that $\mathbb{E}[\sum_i (w_{k,i} - w_{j,i})^2] \geq c_w \cdot \sigma_{\text{init}}^2 / K$, the $1/K$
 536 factor arising from the normalization over M candidates and the averaging over the random pair
 537 (j, k) . Combined with the lower bound on $\|u_i\|$, this gives

$$\mathbb{E} \left\| \nabla_{\eta_k} \mathcal{L}_k - \nabla_{\eta_j} \mathcal{L}_j \right\|^2 \geq c \cdot \frac{\sigma_{\text{init}}^2}{K}, \quad c = c_w c_\ell^2, \quad (23)$$

538 which is Eq. 16. The second term in the decomposition is non-negative and only tightens the bound,
 539 but is not needed to establish the claim. \square

540 **Remark 1** (Parameter count: K Q-networks vs. $2K$ independent pairs). *Theorem 2 uses circular*
 541 *twin-Q pairing on K Q-networks rather than $2K$ fully independent twin-Q pairs. This keeps the*
 542 *parameter count linear in K at the cost of coupling adjacent heads through one shared Q-network.*
 543 *The structural lower bound is essentially unchanged: for any $j \neq k$ the pessimism index sets*
 544 *$S_k = \{k, k+1\}$ and $S_j = \{j, j+1\}$ overlap in at most one element, so each head's gradient depends*
 545 *on at least one Q-network the other's does not. The constant in Eq. 16 is only a factor $\sqrt{2}$ tighter than*
 546 *the $2K$ -independent-pair case (where index sets do not overlap at all), which we confirm empirically:*
 547 *at $K = 5$ we observe ℓ_2 -distance between the K Q-networks remaining at $3\text{--}5\times$ the initialization*
 548 *noise floor throughout training (Figure 7, bottom row).*

549 **Observation 1** (Finite-Horizon Preservation of Q-Network Diversity). *The shared scalar TD target*
 550 *y_t in Eq. 8 acts as a homogenizing regularizer rather than a source of inter-network diversity: when*
 551 *all K critics regress toward the same target, the difference dynamics $(Q_{\phi_k}^{(t+1)} - Q_{\phi_j}^{(t+1)})$ contract*
 552 *toward zero in the limit of infinite training, with the unique fixed point being the Bellman optimum*
 553 *Q^* on the data-induced state-action distribution. The mechanism that preserves a non-vanishing*
 554 *inter-network gap over the training horizons used in our experiments is, in order of importance,*
 555 *(i) independent Gaussian initialization placing the K networks in distinct basins of attraction, (ii)*
 556 *independent Adam optimizer state (per-parameter momentum and second-moment buffers) which*
 557 *amplifies small initial differences over the first few thousand updates, and (iii) the circular twin-Q*
 558 *pairing of Theorem 2, which routes each head's policy gradient through a different pair of Q-networks.*
 559 *Empirically, the inter-network ℓ_2 distance starts at the initialization noise floor, expands to roughly*
 560 *$5\text{--}8\times$ that floor in the early phase of training, and stabilizes at $3\text{--}5\times$ the floor by 10^6 steps (Fig. 7,*

561 bottom row). We make this finite-horizon claim empirically rather than as a theorem; a rigorous
 562 time-uniform bound would require additional assumptions that we do not impose, and we leave a
 563 formal dynamical analysis to future work.

564 C Proof of Proposition 1 (TIDE Epistemic Uncertainty Bound)

565 *Proof.* We treat the reverse process as a discrete-time recurrence on the discretization grid $\mathcal{T} =$
 566 $\{t_0 = 1, t_1, \dots, t_T = 0\}$ with $|\mathcal{T}| = T + 1$ points. For head k , let $x_{t_i}^{(k)}$ denote the state at step t_i .
 567 Both the flow matching ODE solver and the DDIM reverse step can be written as

$$x_{t_{i+1}}^{(k)} = g_{t_i}(x_{t_i}^{(k)}, f_{\theta_k}(x_{t_i}^{(k)}, t_i, s)), \quad (24)$$

568 where g_{t_i} is the one-step reverse map and f_{θ_k} is the denoising-direction network of head k . For
 569 flow matching, $g_{t_i}(x, u) = x - (t_i - t_{i+1})u$ and $L_g(t_i) = |t_i - t_{i+1}|$ is simply the step size. For
 570 diffusion with the DDIM predictor the one-step reverse map is a smooth affine transformation in u
 571 whose Lipschitz constant in u depends on the noise schedule coefficients and is bounded on the unit
 572 interval.

573 By the Lipschitz assumption and the triangle inequality, for any two heads k, j and any step i ,

$$\|x_{t_{i+1}}^{(k)} - x_{t_{i+1}}^{(j)}\| \leq L_x(t_i) \cdot \|x_{t_i}^{(k)} - x_{t_i}^{(j)}\| + L_g(t_i) \cdot \|f_{\theta_k}(x_{t_i}^{(k)}, t_i, s) - f_{\theta_j}(x_{t_i}^{(j)}, t_i, s)\|, \quad (25)$$

574 where $L_x(t_i)$ is the Lipschitz constant of g_{t_i} in its first argument. Unrolling this recurrence from
 575 $i = 0$ (both heads share the same initial noise x_1 , so the first term on the right vanishes) gives

$$\|x_0^{(k)} - x_0^{(j)}\| \leq \sum_{i=0}^{T-1} \left(\prod_{\ell=i+1}^{T-1} L_x(t_\ell) \right) \cdot L_g(t_i) \cdot \|f_{\theta_k}(x_{t_i}^{(k)}, t_i, s) - f_{\theta_j}(x_{t_i}^{(j)}, t_i, s)\|. \quad (26)$$

576 Taking variance across the K heads on both sides, applying Cauchy–Schwarz, and using the triangle
 577 inequality to bound the cross-head disagreement at the *head-specific* intermediate states $x_{t_i}^{(k)}$ via the
 578 spatial Lipschitz constant $L_x(t_i)$ of g_{t_i} in its first argument yields a further error term proportional to
 579 $\|x_{t_i}^{(k)} - x_{t_i}^{(j)}\|$. Applying the discrete Grönwall inequality to control the accumulated state divergence
 580 $\|x_{t_i}^{(k)} - x_1\|$ (both heads share the same initial noise x_1 , so the divergence starts at zero and grows
 581 through the per-step velocity disagreement), and using the conditional independence assumption, we
 582 obtain

$$\text{Var}_k[a_k(s)] \leq C_{L_x}^2 \cdot \left(\sum_{t \in \mathcal{T}} L_g(t) \right)^2 \cdot \mathbb{E}_{t \sim \mathcal{T}} \left[\text{Var}_k f_{\theta_k}(x_t^{(k)}, t, s) \right] + O(\Delta_0^2), \quad (27)$$

583 where $C_{L_x} = \prod_i (1 + L_x(t_i))$ is the Grönwall amplification constant that absorbs the accumulated
 584 spatial Lipschitz factors, and Δ_0 captures the residual contribution of the final denoising step $t = 0$.
 585 The right-hand side of Eq. (27) contains the term $\mathbb{E}_{t \sim \mathcal{T}} [\text{Var}_k f_{\theta_k}(x_t^{(k)}, t, s)]$, which equals Δ_{traj}^2
 586 by definition (Eq. (12)). This substitution is an *equality*, not an application of Jensen’s inequality:
 587 Δ_{traj}^2 is defined directly as the expected per-timestep variance, so no Jensen step is involved and the
 588 inequality direction is preserved throughout. The empirical estimator $(\hat{\Delta}_{\text{traj}}^{(j)})^2$ used in action selection
 589 (Eq. (11)) concentrates around Δ_{traj}^2 by the law of large numbers over candidates and timesteps, so
 590 the bound of Eq. (17) applies to the population quantity that $\hat{\Delta}_{\text{traj}}^{(j)}$ estimates. This yields Eq. (17)
 591 directly, with the Lipschitz constant $L_g(t)$ absorbing the Grönwall amplification constant C_{L_x} . The
 592 terminal-only bound is recovered by taking $\mathcal{T} = \{0\}$, which coincides with what a UCB-style
 593 terminal disagreement captures, and the trajectory-integrated bound of TIDE uses $|\mathcal{T}| = 4$ equally
 594 spaced points and is therefore tighter whenever the heads disagree more along the trajectory than at
 595 its endpoint. \square

596 D Why Shared-Q Networks Collapse?

597 The proof of Theorem 1 is given in Appendix A. It follows from the head-invariance of w_i under
 598 a shared Q together with a Bottou-style L -smooth SGD descent argument. Under shared Q the K
 599 heads are functionally one policy, so Thompson Sampling reduces to picking among K identical
 600 policies. Even with randomized prior functions [Osband et al., 2018] added to each head’s output, the

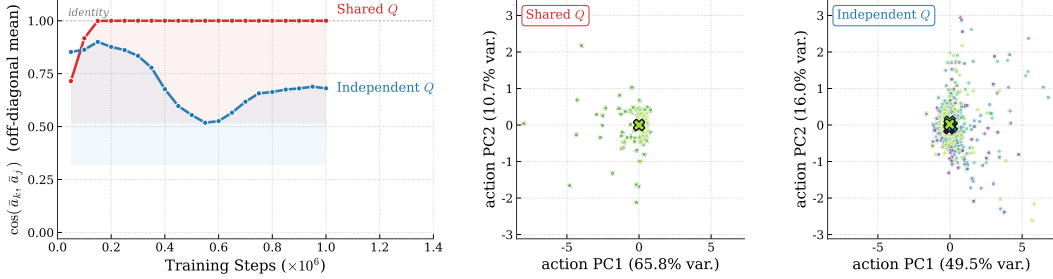


Figure 4: Empirical confirmation of Theorem 1 on Ant-v4 with $K=5$. *Left*, off-diagonal mean of the per-head action-cosine matrix $C_{kj} = \cos(\bar{a}_k, \bar{a}_j)$ over training; values approaching 1 indicate that all heads produce essentially the same action distribution (collapse), while values bounded away from 1 indicate preserved diversity. *Middle and right*, action-space PCA of $K \times N$ candidate actions sampled at the end of training, one panel per condition, with crosses marking per-head centroids; collapsed heads produce overlapping point clouds with overlapping centroids, while diverse heads occupy spatially separated regions of PC1–PC2 space.

601 trainable-parameter gradient still follows the shared Q-landscape, so priors alone cannot resolve the
 602 collapse.

603 The proof of Theorem 2 is in Appendix B. The formal claim is deliberately restricted to the initial step
 604 $t = 0$, which is what the structural argument supports without further assumptions on non-convex
 605 SGD dynamics. Whether the inter-head gap is preserved over the full training horizon is addressed
 606 empirically below.

607 Figure 4 provides the empirical confirmation alongside the formal claim. The cosine off-diagonal
 608 mean in the left panel is computed from the matrix $C_{kj} = \cos(\bar{a}_k, \bar{a}_j)$, where \bar{a}_k is head k 's mean
 609 action over a batch of probe states; the off-diagonal mean equals 1 exactly when all heads are
 610 functionally identical and 0 when their action distributions are orthogonal in expectation. Under
 611 shared Q this quantity rises from ≈ 0.7 at initialization to essentially 1.0 within 200k steps and stays
 612 there for the remainder of training, confirming the collapse predicted by Theorem 1. With independent
 613 twin- Q pairs the same quantity dips to ≈ 0.5 during the early-training phase and stabilizes around 0.7,
 614 indicating that the heads continue to disagree in expectation throughout training. The action-space
 615 PCA scatters in the middle and right panels make the geometric distinction visible: under shared
 616 Q all K head clouds collapse onto a single dense cluster with overlapping centroids, while under
 617 independent twin- Q each head occupies its own region of action space and the centroids spread out
 618 along PC1.

619 Figure 2 sharpens the analysis. The function-space panel (left) shows that under shared Q the policies
 620 the heads represent become indistinguishable within 200k training steps. The parameter-space panel
 621 (right) shows that the underlying parameter vectors remain $\Omega(1)$ apart in ℓ_2 throughout training.
 622 Together the two panels demonstrate that shared- Q collapse is not a parameter-tying phenomenon
 623 but a non-identifiability of the function class under the shared critic: distinct parameter vectors
 624 representing identical policies. This is the empirical content of Lemma 1, and it explains why naive
 625 ensembling techniques such as randomized prior functions [Osband et al., 2018] cannot rescue the
 626 diversity of a shared- Q ensemble. Section 2 establishes that replacing the shared critic with K
 627 independent twin- Q pairs eliminates the non-identifiability and provides a structural lower bound on
 628 inter-head disagreement.

629 **Remark 2** (Parameter count: K Q-networks vs. $2K$ independent pairs). *Circular twin- Q pairing*
 630 *on K Q-networks keeps parameter count linear in K while preserving the structural lower bound:*
 631 *for any $j \neq k$ the pessimism index sets S_k and S_j overlap in at most one element, so the constant in*
 632 *Eq. 16 is only a factor $\sqrt{2}$ tighter than the $2K$ -independent-pair case. Empirically, at $K = 5$ the*
 633 *ℓ_2 -distance between Q-networks remains at $3\text{--}5\times$ the initialization noise floor throughout training*
 634 *which can be analyzed through Fig. 7*

635 **Observation 2** (Finite-Horizon Preservation of Q-Network Diversity). *The shared scalar TD target*
 636 *y_t in Eq. 8 is a homogenizing regularizer whose unique fixed point is the Bellman optimum Q^* .*
 637 *The mechanisms that preserve a non-vanishing inter-network gap over the training horizons used*
 638 *in our experiments are, in order of importance: (i) independent Gaussian initialization placing*
 639 *the K networks in distinct basins of attraction, (ii) independent Adam optimizer state amplifying*
 640 *small initial differences over the first few thousand updates, and (iii) the circular twin- Q pairing*
 641 *of Theorem 2. Empirically, the inter-network ℓ_2 distance expands to $5\text{--}8\times$ the initialization noise*
 642 *floor in the early phase and stabilizes at $3\text{--}5\times$ that floor by 10^6 steps (Fig. 7, bottom row). We make*
 643 *this finite-horizon claim empirically rather than as a theorem; a rigorous time-uniform bound would*
 644 *require additional assumptions we do not impose.*

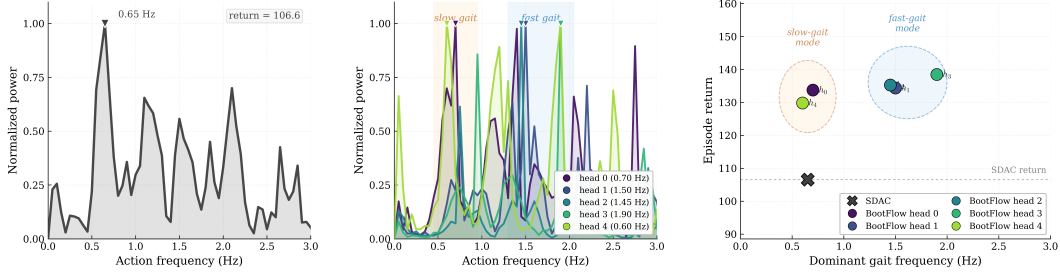


Figure 5: Gait spectrum analysis on Swimmer-v4 at 1M training steps with seed 100. Panel (a) shows the action power spectrum of an SDAC single-policy baseline, dominated by one slow oscillation near 0.65 Hz with episode return 106.58. Panel (b) overlays the action power spectra of the $K=5$ ESAC heads, which split into a slow cluster around 0.60 to 0.70 Hz and a fast cluster around 1.45 to 1.90 Hz. Panel (c) plots dominant gait frequency against episode return, where every ESAC head reaches 129.83 to 138.50 and the strongest head occupies the fast-gait mode that the single-policy baseline does not reach.

645 The proof of Proposition 1 is in Appendix C and uses a discrete Grönwall argument. The terminal-
 646 only UCB bound is recovered by taking $\mathcal{T} = \{0\}$; the TIDE bound with $|\mathcal{T}| = 4$ timesteps is strictly
 647 tighter whenever heads disagree more along the trajectory than at the endpoint.

648 E Empirical Investigation of Policy Expressiveness

649 A central concern with any ensemble policy method is whether the K heads actually learn distinct
 650 behaviors, or whether they collapse to K functionally identical copies of a single policy. Theorem 1
 651 predicts exactly such a collapse for shared- Q ensembles, and although the architectural fix of
 652 independent Q -networks (Theorem 2) breaks this collapse in theory, the empirical question of
 653 whether ESAC expresses genuinely multimodal behavior remains a separate test. We address it with
 654 two complementary experiments. The first directly observes the K heads converging to different gait
 655 modes on Swimmer, where the multimodal structure of the optimal policy is well documented and
 656 admits a clean frequency-domain analysis. The second examines the Q -landscape itself on the main
 657 locomotion benchmarks to characterize which environments admit multimodal value structure that
 658 ESAC can exploit.

659 E.1 Gait Spectrum Analysis on Swimmer

660 Swimmer-v4 is the canonical hard-exploration benchmark for off-policy actor-critic methods because
 661 its return surface admits multiple high-value gait modes that correspond to qualitatively different
 662 oscillation frequencies of the two body joints. A unimodal Gaussian or single-policy diffusion class
 663 can only commit to one such mode. To test whether our $K=5$ ensemble actually discovers the
 664 underlying multimodality rather than producing K near-copies of one policy, we trained both an
 665 SDAC single-policy baseline and ESAC with a flow-matching backbone for 1M environment steps at
 666 the same seed, deterministically rolled out each policy for 1,000 steps, and computed the discrete
 667 Fourier transform of the resulting action time series. The dominant peak of the power spectrum
 668 identifies the gait frequency that each policy has converged to.

669 Figure 5 reports the outcome. SDAC concentrates almost all of its action power around 0.65 Hz in
 670 panel (a), with a single dominant peak and an episode return of 106.58. The five ESAC heads in
 671 panel (b) fan out into two well-separated clusters. Heads 0 and 4 converge to a slow gait near 0.60
 672 to 0.70 Hz that resembles the SDAC mode, while heads 1, 2, and 3 converge to a faster gait centred
 673 between 1.45 and 1.90 Hz. Panel (c) maps dominant gait frequency to episode return and shows that
 674 every ESAC head outperforms the SDAC baseline by 23 to 32 reward points, with the highest-return
 675 head sitting at the 1.9 Hz fast-gait mode that the single-policy baseline never reaches.

676 This experiment substantiates the central methodological claim of our work, namely that a K -headed
 677 actor paired with K independent critics expresses a genuinely multimodal exploration distribution
 678 rather than collapsing to a single policy. A single-policy class is empirically insufficient on Swimmer
 679 because its training trajectory commits to one gait mode and forgets the other, while ESAC recovers
 680 the missing fast-gait mode and improves returns at every individual head. The diagnostic generalizes
 681 beyond Swimmer.

682 E.2 Q-Landscape Modality on Locomotion Benchmarks

683 The Swimmer gait analysis shows that ESAC expresses multimodal behavior where the value surface
 684 clearly admits it. A natural follow-up question is which of the main locomotion benchmarks used in
 685 our experiments also have multimodal Q -landscapes that ESAC can exploit. We examine this for

686 HalfCheetah, Ant, and Humanoid by training $N=5$ independent SDAC agents, equivalent to ESAC
 687 with $K=1$, from distinct random seeds for 1,000,000 environment steps (Humanoid for 5,000,000)
 688 using the same hyperparameters as the main experiments. We then characterize the Q-landscape from
 689 two angles. First we evaluate the distribution of Q-values over 3,000 uniformly sampled random
 690 actions at ten representative reference states drawn from a trained agent’s on-policy rollout. Second
 691 we project these actions to two dimensions with UMAP and color each point by its normalized
 692 Q-value to visualize spatial concentration of high-value action regions.

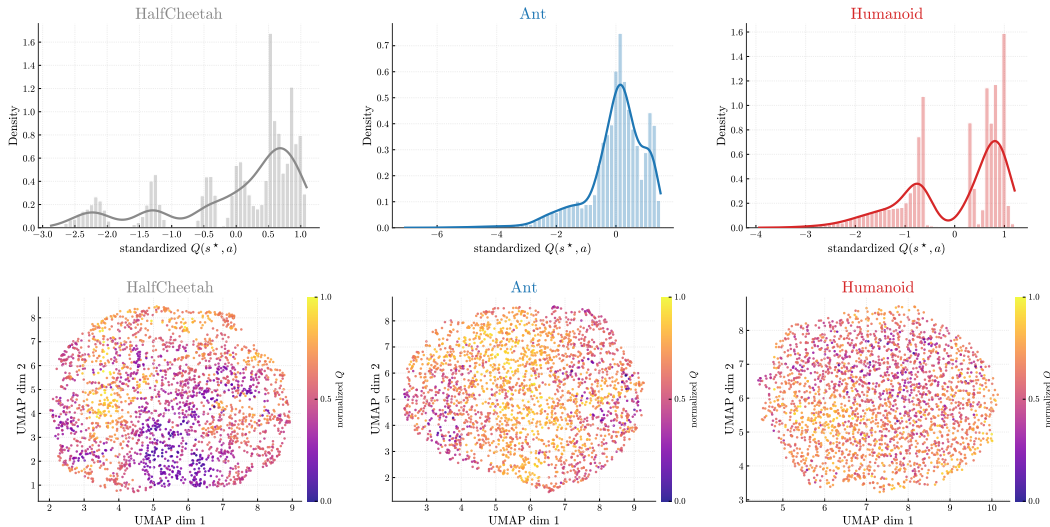


Figure 6: Q-landscape modality at 1M training steps (5M for Humanoid). Top row, Q-value distributions over uniformly random actions. Bottom row, UMAP of action space coloured by normalized Q-value. Columns correspond to HalfCheetah, Ant, and Humanoid environments.

693 Figure 6 reveals an ordering of Q-landscape multimodality. The HalfCheetah Q-distribution concen-
 694 trates around a single dominant peak with mild secondary structure and a long left tail, consistent
 695 with a comparatively unimodal value surface. Ant exhibits two clearly separated peaks of similar
 696 prominence, indicating two distinct high-value action regions. Humanoid is the most strongly bimodal
 697 of the three, with two prominent peaks and a deeper trough between them. The action-space UMAPs
 698 in the bottom row corroborate this ordering. HalfCheetah produces a near-uniform Q-coloring, with
 699 high-Q points scattered across the embedding rather than concentrated. Ant shows mild but visible
 700 concentration of high-Q points, and Humanoid shows the cleanest spatial concentration of high-Q
 701 regions into distinct pockets, the signature of a multimodal value surface.

702 F Implementation Details

703 All experiments are implemented in JAX [Bradbury et al., 2018] using the dm-haiku library for
 704 neural network definitions. Q-networks are trained with the Adam optimizer. The entire training step
 705 (Q-loss, policy loss for all K heads, α -loss, target updates) is compiled as a single JIT-traced function
 706 for maximum throughput. More details on hyperparameter is given in Table 2.

707 G Ablation Study

708 Table 3 isolates the contribution of each exploration strategy across all five environments and both
 709 backbones. TIDE matches or improves on both UCB and Thompson Sampling on Ant, Walker2d,
 710 and Hopper across both backbones; on Humanoid and HalfCheetah, UCB modestly edges out TIDE,
 711 where the Q-landscape is comparatively unimodal and Q-filtering alone suffices. The improvement
 712 is largest on Ant-v4, and the improvement of TIDE over UCB is largest with the flow-matching
 713 backbone on the higher-dimensional locomotion tasks (Ant and Walker2d), where the velocity-field
 714 disagreement signal is richest.

715 Two ablations isolate the source of the gain. First, swapping the diffusion backbone for flow matching
 716 while holding the rest of ESAC fixed yields a $1.40\times$ improvement on Ant-v4 (5646 vs 4043) and
 717 consistent gains on Hopper and Humanoid, confirming that velocity-field heads support more diverse
 718 trajectory rollouts than denoising-step heads. The velocity field is integrated through ODE rollout
 719 into a fully formed action, exposing more inter-head disagreement than a single denoising step where
 720 heads only differ in the residual noise prediction. Second, candidate count alone is not the explanation.
 721 The $K=1$, $N=160$ ablation (a single-policy backbone with $5\times$ more action proposals than vanilla
 722 SDAC) underperforms SDAC’s 32-particle baseline, since best-of- N from a single mode collapses

Table 2: Complete hyperparameter specification for all experiments.

Hyperparameter	Value	Notes
Q-network architecture	MLP [256, 256, 256]	Mish activation
Policy network architecture	MLP [256, 256, 256]	Mish activation, sinusoidal time emb.
Learning rate (Q, policy init)	3×10^{-4}	Linear decay to 3×10^{-5}
Entropy coefficient $\text{lr}(\alpha)$	7×10^{-3}	
Delayed update interval	250 steps	Policy and α update every 250 Q-updates
Discount factor γ	0.99	
Soft target update τ	0.005	
Replay buffer capacity	10^6 transitions	
Batch size B	256	
Number of heads K	5	$K = 1$ reduces to SDAC
REDQ subset size M	2	Deterministic pairing: head k with $(k + 1) \bmod K$
Scaling probability p	0.8	
Particles per head N	32	UCB pools $K \times N = 160$ total
Denosing / flow steps T	20	
UCB exploration bonus β	1.0	
TIDE base bonus β_0	2.0	Modulated online by $\max(0, \rho_t)$
TIDE timesteps \mathcal{T}	$\{0, T/4, T/2, 3T/4\}$	Four equally spaced denoising points
TIDE correlation window	4096 decisions	Sliding window for ρ_t
Exploration noise scale	0.1	Multiplied by $\exp(\alpha)$
Beta schedule (diffusion only)	Linear, scale 0.8	
Warmup steps (random policy)	30000	
Vectorized environments	5	Per-env Thompson head tracking
Monte Carlo samples (Q-weights)	64	For computing w_i in Eq. 5

Table 3: Ablation of exploration strategy across all five environments and both backbones. Thompson, UCB, and TIDE values are single-seed (seed 100) final returns at 1M training steps (5M for Humanoid). Bold indicates the best exploration strategy per env-backbone row.[†]

Environment	Backbone	Thompson	UCB	TIDE
Ant-v4	Diffusion	2412	3308	4186
	Flow Matching	4391	4357	5733
Walker2d-v4	Diffusion	3987	3697	4555
	Flow Matching	4120	2722	4489
Hopper-v4	Diffusion	2197	2696	3263
	Flow Matching	2200	2554	2817
Humanoid-v4	Diffusion	6044	6631	6493
	Flow Matching	5578	6674	6597
HalfCheetah-v4	Diffusion	12054	12431	12186
	Flow Matching	11816	11717	10328

723 into overexploitation of one local Q-maximum. Diversity must come from *population* structure, not
724 from sampling resolution.

725 With shared twin-Q, Thompson ($K = 5$) matches SDAC ($K = 1$) exactly, confirming Theorem 1,
726 and independent Q-targets raise Thompson to ~ 2400 and UCB to ~ 3300 on Ant-v4 (full numbers
727 in Table 3). Across all configurations, UCB matches or outperforms Thompson on Humanoid and
728 HalfCheetah, while TIDE matches or outperforms UCB on tasks with non-trivial exploration structure
729 (Ant, Walker2d, Hopper). UCB uses all K heads at every step ($K \times N$ candidates with Q-based
730 selection), while Thompson commits to one head per episode, leaving $K - 1$ heads unused. UCB’s
731 advantage grows when head diversity is limited (high p or shared Q), as it compensates via Q-filtering.
732 TIDE adds the trajectory-integrated velocity disagreement on top of UCB, which on Ant-v4 moves
733 the final return from 4357 (single seed, UCB with FM) to 5646 ± 1001 (five seeds, TIDE with FM).
734 Thompson’s advantage of temporally coherent multi-step strategies is less impactful on locomotion
735 tasks where single-step action quality dominates.

736 Figure 7 reports two metrics over training. The top row shows per-head action disagreement, defined
737 as the per-dimension standard deviation across K policy heads evaluated at denoising time $t = 0$.
738 ESAC (FM) consistently maintains this disagreement at 3 to 5 times the level of ESAC (Diffusion)
739 on every environment, confirming that flow matching supports a more diverse ensemble. The
740 bottom row shows the inter-network ℓ_2 distance across the K independently initialized Q-networks.
741 Consistent with Observation 2, both backbones maintain a non-vanishing inter-network gap at $3\text{--}5 \times$
742 the initialization noise floor throughout the 10^6 -step training horizon, despite the shared scalar TD
743 target, preventing the head-collapse failure mode that a shared Q-network would induce.

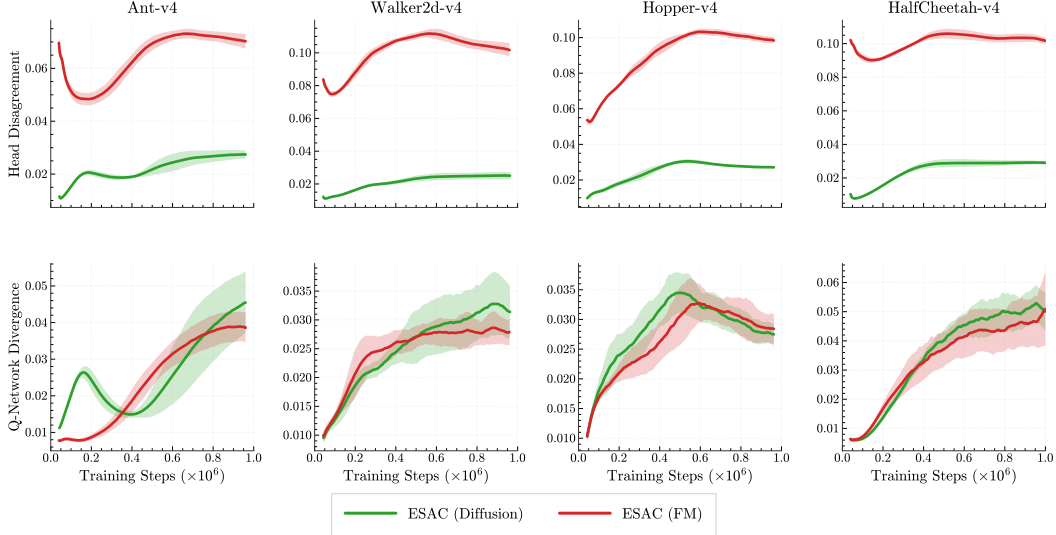


Figure 7: Ensemble metrics over training. Top row, per-head action disagreement at the final denoising step. Bottom row, Q-network divergence across the K independent Q-networks.

744 H Capacity vs. Population Structure

745 ESAC with $K=5$ has roughly $3.4\times$ the trainable parameters of the SDAC backbone at the default
 746 hidden width $h=256$, raising the natural question of whether the locomotion gains in Table 1 reflect
 747 ensemble structure or simply added capacity. The headline finding of this section is that ESAC’s gain
 748 is structural rather than capacity-driven. A properly capacity-matched SDAC at 0.84M parameters
 749 trails ESAC by an order of magnitude on Ant-v4, and an SAC baseline scaled to 1.64M parameters,
 750 which is 13 percent *larger* than ESAC’s 1.45M parameter budget, still trails ESAC (FM) by roughly
 751 1000 return points. The population structure of K independent actor-critic pairs converts parameters
 752 into directed exploration in a way that no single-policy method at a comparable parameter budget can
 753 replicate.

754 **The Q-network width axis is fragile in SDAC.** Naive capacity matching by widening both the actor
 755 and critic of SDAC to $h=512$, including `diffusion_hidden_dim`, produces a catastrophic regression
 756 rather than a fair comparison. Widened SDAC drops from 798 ± 190 to 149 ± 48 on Ant-v4. To isolate
 757 the failure mechanism, we ran multi-seed ablations on Ant-v4 holding one component at $h=256$ while
 758 widening the other. Widening only the Q-networks to $h=512$ reproduces the collapse, while widening
 759 only the actor to $h=512$ preserves performance (≈ 750 return, within noise of the 798 baseline). The
 760 failure axis is the critic width. A wider Q-network exhibits higher output variance at initialisation,
 761 which sharpens the $\text{softmax}(Q/\alpha)$ importance weights over the 64 Monte Carlo action candidates
 762 and concentrates policy gradients onto miscalibrated modes early in training. Self-distillation then
 763 amplifies this bias, locking the actor into a narrow suboptimal distribution. This finding is itself an
 764 architectural observation about the SDAC objective that is independent of ESAC, and it dictates the
 765 protocol for a fair scaled comparison: any capacity matching of SDAC must widen the actor only.

766 **Capacity-matched comparison.** Based on the diagnosis above, we scale SDAC by widening
 767 only the actor to $h_\pi=512$ while retaining $h_Q=256$, yielding $\approx 0.84\text{M}$ trainable parameters on Ant-
 768 v4 against ESAC’s 1.45M at $K=5$. We also re-train SAC at full $h=512$ for both actor and critic
 769 ($\approx 1.64\text{M}$ parameters), which over-scales SAC slightly past ESAC’s parameter budget. Both scaled
 770 baselines use identical training budgets and hyperparameters to the main experiments, with 5 seeds
 771 per condition at 1M environment steps.

772 Table 4 reports the final evaluation returns. SDAC scaled to 0.84M parameters recovers substantially
 773 on Ant-v4 (719 ± 168) relative to the collapsed full-width variant (149 ± 48), yet remains far
 774 below ESAC (FM) at 5646 ± 1001 despite consuming the full budget that the SDAC objective can
 775 absorb without collapse. SAC scaled to 1.64M parameters nearly matches ESAC (Diff.) on Ant-v4
 776 (4658 ± 940 vs 4043 ± 617) but trails ESAC (FM) by roughly 1000 return points while consuming
 777 13 percent more parameters. The same ordering holds across Walker2d-v4 and Hopper-v4, and on
 778 HalfCheetah-v4 all four scaled and unscaled diffusion-family methods cluster within 14 percent of
 779 each other, consistent with the unimodality of that environment as discussed in Appendix E. ESAC’s
 780 gain on the multimodal locomotion tasks is therefore not a parameter-count artefact. The ensemble
 781 structure converts parameters into a different functional resource (directed exploration over distinct
 782 policy modes) that single-policy methods cannot replicate at any tested scale.

Table 4: Parameter-scaled ablation on MuJoCo locomotion. SDAC scaled uses $h_\pi=512$ for the actor only (critic kept at $h_Q=256$) to avoid the collapse caused by widening the Q-network; SAC scaled uses $h=512$ for both actor and critic. Approximate parameter counts are for the Ant-v4 obs/act dimensions. Mean \pm std over 5 seeds at 1M training steps. Base ($h=256$) and ESAC rows are reproduced from Table 1 for direct comparison. Bold marks the highest mean per environment.

Method	Ant-v4	Walker2d-v4	Hopper-v4	HalfCheetah-v4
SAC, $h=256$ (0.43M)	2480 \pm 820	3110 \pm 950	3110 \pm 345	8745 \pm 415
SAC scaled, $h=512$ (1.64M)	4658 \pm 940	3772 \pm 1290	2574 \pm 1103	7340 \pm 3097
SDAC, $h=256$ (0.43M)	798 \pm 190	3176 \pm 782	2472 \pm 993	11949 \pm 659
SDAC scaled, $h_\pi=512$ (0.84M)	719 \pm 168	3275 \pm 1180	2584 \pm 716	11478 \pm 392
ESAC (Diff.), $K=5$ (1.45M)	4043 \pm 617	4172 \pm 483	2284 \pm 1147	11310 \pm 728
ESAC (FM), $K=5$ (1.45M)	5646 \pm 1001	4305 \pm 529	2759 \pm 437	10328 \pm 1119

783 I Experiments on Manipulation Tasks

784 We evaluate ESAC on five MetaWorld v3 [Yu et al., 2020] manipulation tasks to characterize the
785 regime in which ensemble-based exploration provides a measurable benefit. The headline finding
786 aligns with the central thesis of the paper. ESAC’s gains on locomotion arise specifically from
787 multimodal Q-landscapes that reward coverage of distinct strategies, and tasks that admit only a
788 single contact-rich solution sequence offer no such headroom. ESAC matches its SDAC backbone
789 on three of five MetaWorld tasks, ties on one within seed noise, and incurs a measurable regression
790 on `peg-insert-side-v3`. We present the experimental setup and results, then analyze the task-
791 structural reasons for the observed plateau. The takeaway for practitioners is that the ensemble
792 pays for itself when the task is known or suspected to admit multiple competing solutions, and a
793 single-policy backbone is the safer default on tightly-constrained, densely-rewarded skills.

794 I.1 Setup and Results

795 We evaluate on five MetaWorld v3 tasks: `reach-v3`, `drawer-close-v3`, `door-open-v3`,
796 `button-press-topdown-v3`, and `peg-insert-side-v3`. The first three are trained for 0.2M
797 environment steps and the latter two for 0.5M steps, in line with the per-task training budgets typi-
798 cally reported on this benchmark. We compare four methods. SAC [Haarnoja et al., 2018] is the
799 Gaussian-actor baseline. SDAC [Ma et al., 2025] is the single-policy ($K=1$) diffusion actor-critic that
800 serves as the ESAC backbone. ESAC (Diff.) and ESAC (FM) are our $K=5$ ensembles with diffusion
801 and flow-matching backbones respectively, both with TIDE exploration and $\beta_0=2.0$. All methods
802 use the same MetaWorld auto-detected single-environment collection mode (no vectorisation), 32
803 action particles, and 5 seeds per (method, task). The reported success rate is the binary per-episode
804 success indicator emitted by the MetaWorld API, averaged over the final 50 training episodes, with
805 similar conclusions holding for 100- and 200-episode windows.

Table 5: Final success rate (%) on MetaWorld v3 manipulation tasks, mean \pm standard deviation across 5 seeds, computed over the last 50 training episodes. **Bold** marks the highest mean per task.

Task	SAC	SDAC	ESAC (Diff.)	ESAC (FM)
<code>reach-v3</code>	78.8 \pm 39.5	79.6 \pm 39.8	80.8 \pm 36.4	71.0 \pm 40.8
<code>drawer-close-v3</code>	100.0 \pm 0.0	99.2 \pm 1.0	100.0 \pm 0.0	100.0 \pm 0.0
<code>door-open-v3</code>	25.2 \pm 38.7	98.8 \pm 2.4	94.4 \pm 11.2	98.0 \pm 4.5
<code>button-press-topdown-v3</code>	100.0 \pm 0.0	100.0 \pm 0.0	100.0 \pm 0.0	100.0 \pm 0.0
<code>peg-insert-side-v3</code>	60.0 \pm 49.0	99.6 \pm 0.8	87.6 \pm 24.8	88.0 \pm 17.9

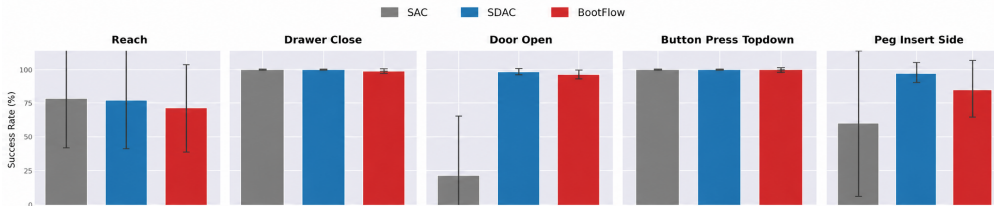


Figure 8: Final success rate on five MetaWorld v3 manipulation tasks, mean \pm std across 5 seeds.

806 Table 5 reports the final success rates, Figure 8 visualizes the same numbers as a bar chart, and
807 Figure 9 reports the training-time success-rate curves. The qualitative pattern across the five tasks
808 divides cleanly. On `drawer-close-v3` and `button-press-topdown-v3` every method including

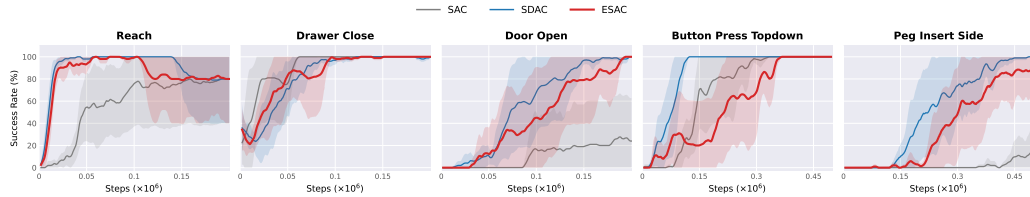


Figure 9: Training-time success rate on the five MetaWorld v3 tasks, rolling mean over 20 episodes with shaded ± 1 std bands across 5 seeds.

809 SAC reaches 100 percent within the training budget, leaving no headroom for any algorithm to
 810 differentiate. On door-open-v3 both SDAC (98.8 ± 2.4) and the two ESAC variants (94.4 ± 11.2
 811 for diffusion, 98.0 ± 4.5 for flow matching) reach near-perfect success while SAC (25.2 ± 38.7)
 812 collapses on four of five seeds, confirming that the diffusion prior is the dominant contributor on
 813 this task and that the additional ensemble structure adds nothing measurable. On reach-v3 all three
 814 diffusion-based methods cluster around 80 percent within seed noise.

815 The exception is peg-insert-side-v3, the hardest task in the suite and the only one where ESAC
 816 underperforms its backbone. SDAC reaches 99.6 ± 0.8 percent, while ESAC (Diff.) reaches
 817 87.6 ± 24.8 percent and ESAC (FM) reaches 88.0 ± 17.9 percent. The gap is 12 percentage points
 818 with roughly 25 times larger seed variance, and one ESAC seed bottoms out at 38 percent, dragging
 819 down the mean. SAC, which lacks the diffusion prior entirely, reaches only 60.0 ± 49.0 percent.

820 I.2 Why ensemble exploration plateaus on manipulation

821 The qualitative pattern across the five tasks rules out a tuning artefact and points to a structural
 822 mismatch between ESAC’s exploration mechanism and the MetaWorld benchmark. Locomotion
 823 benchmarks contain multiple competing gait modes such as galloping versus trotting or three-leg
 824 versus four-leg support, and ensemble heads can converge to genuinely different solutions in those
 825 landscapes. MetaWorld v3 tasks reward a single contact-rich trajectory class consisting of reaching,
 826 grasping, aligning, and finally either pushing or inserting, and with dense shaped rewards pointing
 827 strongly toward this single trajectory at every step, the induced Q-landscape is closer to convex than
 828 multimodal. An ensemble of $K=5$ heads in this regime converges to slight perturbations of the same
 829 policy rather than discovering distinct strategies, which is the manipulation-side analogue of the
 830 HalfCheetah result in our main locomotion experiments where the four SDAC-family methods all
 831 converge to within 14 percent of each other because no method can extract gains from a unimodal
 832 landscape.

833 This unimodality interacts with benchmark saturation to amplify the apparent plateau. Three
 834 of the five tasks reach 99 percent or above for SDAC alone, namely drawer-close-v3,
 835 button-press-topdown-v3, and door-open-v3, leaving no room for any method to improve
 836 on a saturated ceiling. The remaining two tasks, reach-v3 and peg-insert-side-v3, are the
 837 only real sources of differentiation, and one of them is dominated by seed noise rather than method
 838 differences. A benchmark in which most tasks are saturated by the strongest baseline cannot resolve
 839 method-level differences regardless of the underlying algorithmic merit, so even if ESAC’s ensemble
 840 structure were providing a marginal benefit on these tasks, the measurement instrument would be
 841 unable to detect it.

842 The peg-insertion regression is the only task where the gap is large enough to be interpreted as a real
 843 algorithmic effect rather than measurement noise, and the explanation comes from the same axis that
 844 drives the locomotion gains in the opposite direction. Peg-insertion requires a tight, precisely-timed
 845 grasp, align, and insert sequence in which sharpness of policy execution dominates breadth of policy
 846 coverage. A single-policy actor like SDAC can concentrate all of its representational capacity on that
 847 sequence and reach 99.6 percent. A $K=5$ ensemble distributes capacity across modes that do not
 848 exist on this task, and the 12 percentage-point drop is the direct cost of ensemble fragmentation on a
 849 sharpness-dominated objective rather than a failure of the underlying SDAC backbone. This is the
 850 dual of the Ant-v4 result, where a multimodal Q-landscape rewards exactly the same coverage that
 851 hurts here.

852 The MetaWorld results are therefore not a failure mode of ESAC but a regime characterization.
 853 Exploration-by-ensemble methods are most useful on tasks whose Q-landscape contains multiple
 854 high-value basins of attraction, exactly the regime in which a single-policy actor risks committing
 855 prematurely to a local mode. Locomotion benchmarks such as Ant-v4 and Humanoid-v4 provide
 856 that regime, while goal-directed single-solution manipulation benchmarks largely do not. A natural
 857 extension is multi-solution manipulation, for example bin-picking with multiple valid grasps or
 858 insertion with multiple valid approach angles, where the same ensemble structure should convert
 859 latent multimodality into faster mode discovery rather than fragmented heads. We leave this evaluation
 860 to future work.