

---

# LASER: Latent Space Adjoint Matching for Support Constrained Entropy Regularized Offline RL

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 While offline reinforcement learning (RL) enables policy optimization from static  
2 datasets without costly online interaction, it remains bottlenecked by the risk of  
3 executing out-of-distribution (OOD) actions. Recent approaches mitigate this by  
4 learning a flow matching behavior cloning policy to approximate the dataset, and  
5 then performing RL within its constrained latent space. However, naïvely doing  
6 gradient descent in the latent space can easily cause the policy to collapse into  
7 a brittle mode or exploit sharp artifacts of the learned critic. In this work, we  
8 identify that entropy regularization is essential in latent-space RL for addressing  
9 these challenges. We introduce LASER, a novel offline RL algorithm that applies  
10 latent-space adjoint matching to achieve entropy-regularized latent space RL with  
11 expressive flow policies without backpropagation through time. Through com-  
12 prehensive experiments on 40 challenging OGBench tasks with varying dataset  
13 qualities, we show that LASER achieves state-of-the-art performance. Notably,  
14 LASER uses a single, *constant* set of hyperparameters to outperform baselines even  
15 with their hand-tuned hyperparameters, highlighting its robust applicability.

## 16 1 Introduction

17 Offline reinforcement learning (RL) aims to learn policies from a static dataset that improve upon the  
18 behavior policy that generated the dataset without any online interaction [1]. This creates a central  
19 tension: improving upon the behavior policy requires the learned policy to deviate from it, but still  
20 output actions where the dataset provides reliable supervision. For a fixed state, actions outside  
21 the dataset support may receive spuriously high predicted values, which can be exploited by policy  
22 optimization to yield out-of-distribution (OOD) actions that perform poorly when deployed [1].

23 A major line of work addresses the problem of OOD actions by constraining policy improvement to  
24 remain close to the behavior distribution, by regularizing a statistical distance to the behavior policy  
25 [2, 3, 4, 5, 6]. However, depending on the choice of statistical distance, these methods can be either  
26 too restrictive and produce an overly conservative policy (e.g., KL divergence), or not restrictive  
27 enough and still produce OOD actions (e.g., Wasserstein distance, MMD distance) [7].

28 This has motivated a more direct alternative: enforce **support constraints** on the policy extraction  
29 process to guarantee that the actions produced lie within the support of the action distribution of  
30 the behavior policy [8, 7]. In particular, these methods learn a flow-based generative model using  
31 flow-matching [9, 10, 11] to imitate the behavior policy, resulting in a compact latent space that maps  
32 to the behavior policy. Assuming a flow-based model that perfectly captures the behavior policy, any  
33 sample from the latent space is a sample from the behavior policy and thus lies within the support of  
34 the dataset. Policy optimization performed on this latent space, either via gradient descent (as in [7])  
35 or SAC [12] (as in [8]) yields a policy that satisfies support constraints by construction.

36 In this work, we argue that support constraints are only part of the story. Once policy optimization is  
 37 restricted to the support, the central difficulty becomes policy extraction: how should one reliably  
 38 select high-value actions from the support without collapsing to brittle modes or exploiting sharp  
 39 artifacts of the learned critic? Even when actions remain in support, the critic may be sharp, irregular,  
 40 or locally overestimated, and a low-entropy extraction procedure can collapse onto brittle modes.

41 We identify **entropy regularization** as a missing principle in support-constrained flow policy op-  
 42 timization. While entropy regularization has been thoroughly studied in online RL for its role in  
 43 improving exploration [13, 12] and robustness [13, 14], its role in offline RL, and in particular in  
 44 support-constrained policy extraction, has not been isolated and made explicit. We argue that entropy  
 45 regularization similarly plays a crucial role in offline RL, not only in discouraging premature collapse  
 46 of the policy to local minima, but also in smoothing the policy optimization landscape by regularizing  
 47 the critic [15].

48 To this end, we propose **Latent Space Adjoint Matching for Support Constrained Entropy Regularized**  
 49 **Offline RL (LASER)**, a framework for entropy-regularized support-constrained policy extraction with  
 50 latent-space flow policies. LASER builds on the latent support parameterization used in recent flow-  
 51 based offline RL methods [7, 8], but replaces implicit or architecture-dependent regularization with  
 52 an explicit entropy-regularized policy improvement objective. The main technical challenge is that  
 53 computing and optimizing the entropy of a flow policy in latent space is nontrivial. We address this  
 54 using **adjoint matching** [16], which enables entropy-regularized optimization of expressive flow  
 55 policies *in latent space* while preserving the support constraints induced by the learned latent space.

56 We make the following contributions:

- 57 1. We identify entropy regularization as a central ingredient for stable support-constrained policy  
 58 extraction with flow-matching policies in offline RL.
- 59 2. We reinterpret recent support-constrained flow policy methods, including DSRL and ReFORM,  
 60 through the lens of implicit entropy regularization.
- 61 3. We introduce **LASER**, a principled entropy-regularized objective for flow-based offline policy  
 62 improvement that preserves dataset-support constraints while enabling smoother and more reliable  
 63 policy extraction.
- 64 4. Extensive experiments demonstrate that LASER, *using a fixed set of hyperparameters*, dominates all  
 65 flow policy-based baselines using their per-task hand-tuned hyperparameters on the performance  
 66 profile curve.

## 67 2 Preliminaries

### 68 2.1 Offline RL

69 We use  $\Delta(\mathcal{X})$  to denote the set of probability distributions over  $\mathcal{X}$ , and grey notations for place-  
 70 holder variables. A Markov Decision Process (MDP) is defined by  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, \rho_0, P, \gamma)$ ,  
 71 with state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , reward function  $r(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , initial state dis-  
 72 tribution  $\rho_0 \in \Delta(\mathcal{S})$ , transition dynamics function  $P(s'|s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ , and dis-  
 73 count factor  $\gamma \in [0, 1]$ . Given a dataset of trajectories  $\mathcal{D}$  generated by some behavior policy  
 74  $\pi_\beta(a|s) : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ , offline RL aims to learn a policy  $\pi_A(a|s) : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ , that maxi-  
 75 mizes the expected discounted return  $R(\pi_A) = \mathbb{E}_{\tau \sim \rho^{\pi_A}(\tau)} \left[ \sum_{h=0}^H \gamma^h r(s_h, a_h) \right]$ , where  $\rho^{\pi_A}(\tau) =$   
 76  $\rho_0(s_0) \pi_A(a_0|s_0) P(s_1|s_0, a_0) \cdots \pi_A(a_H|s_H)$ . Note that sampling trajectories during training is not  
 77 allowed in the offline RL setting.

78 **The OOD Problem.** A primary challenge in offline RL is OOD actions [1]. Standard actor-  
 79 critic frameworks evaluate a policy  $\pi_A$  by estimating its state-action value function ( $Q$ -function)  
 80  $Q(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , representing the expected discounted return:

$$Q^{\pi_A}(s, a) = \mathbb{E} \left[ \sum_{h=0}^H \gamma^h r(s_h, a_h) \mid s_0 = s, a_0 = a, a_h \sim \pi_A(s_h), \forall h \geq 1 \right]. \quad (1)$$

81 To learn  $Q^{\pi_A}$ , fitted Q-evaluation is typically used to minimize the temporal difference (TD) error:

$$L(\phi) = \mathbb{E}_{(s, a, s') \sim \mathcal{D}, a' \sim \pi_A(s')} \left[ \left( r(s, a) + \gamma Q_\phi^{\pi_A}(s', a') - Q_\phi^{\pi_A}(s, a) \right)^2 \right], \quad (2)$$

82 where  $Q_\phi^{\pi_A}$  is a target network often updated via Polyak averaging [17]. However, because the dataset  
 83  $\mathcal{D}$  is fixed, querying the target network with unfamiliar actions  $a' \sim \pi_A(s')$  can yield unreliable, and  
 84 potentially large  $Q$ -values. Consequently, the policy may erroneously optimize toward these OOD  
 85 actions [1]. To address this issue, the learned policy should satisfy the **support constraint** [5]:

$$\text{supp}(\pi_A(\cdot|s)) \subseteq \text{supp}(\pi_\beta(\cdot|s)), \quad \forall s, \quad (3)$$

86 where  $\text{supp}(p) := \{x \mid p(x) > 0\}$  denotes the support of a distribution  $p$ . To achieve this, many prior  
 87 offline RL algorithms regularize the statistical distance  $D(\pi_A \parallel \pi_\beta)$  between the learned policy  $\pi_A$   
 88 and the behavior policy  $\pi_\beta$ , resulting in the following objective for policy updates:

$$\min_{\pi_A} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_A(s)} \left[ -Q_\phi^{\pi_A}(s, a) + \alpha D(\pi_A \parallel \pi_\beta) \right], \quad (4)$$

89 where  $\alpha$  is a hyperparameter for the regularization weight. However, these approaches may not fully  
 90 avoid OOD or can limit the improvement of  $\pi_A$  w.r.t.  $\pi_\beta$  depending on the choice of the statistical  
 91 distance  $D$  [7], and the weight  $\alpha$  requires hand-tuning for different tasks and datasets [1, 18].

## 92 2.2 Mitigating OOD Actions with Latent Space RL

93 To mitigate the limitations of statistical distance regularization, recent works propose learning policies  
 94 in a *latent space*  $\mathcal{L}$  [19, 8, 7]. Specifically, a **decoder**  $\mu_{\theta_{L \rightarrow A}}^s(z) : \mathcal{L} \times \mathcal{S} \rightarrow \mathcal{A}$  is learned that maps  
 95 samples  $z \sim q_L$  drawn from the latent space prior distribution  $q_L(z) \in \Delta(\mathcal{L})$  to the dataset distribution  
 96  $\pi_\beta(\cdot|s)$ , i.e., the pushforward of the prior  $q_L$  through the decoder  $\mu_{\theta_{L \rightarrow A}}$  approximates the behavior  
 97 policy  $\pi_\beta(\cdot|s)$ . Then, instead of learning a policy  $\pi_A(\cdot|s)$  in the action space  $\mathcal{A}$ , we can learn a  
 98 **latent-space policy**  $\pi_L(z|s) : \mathcal{S} \rightarrow \Delta(\mathcal{L})$  which induces a full policy in the action space via the  
 99 pushforward through the decoder  $\mu_{\theta_{L \rightarrow A}}$ , i.e.,

$$\pi_A(\cdot|s) := (\mu_{\theta_{L \rightarrow A}}^s)_\# \pi_L(\cdot|s), \quad (5)$$

100 where  $(\cdot)_\#$  is the pushforward operator. This parametrization (5) transforms the difficult action-space  
 101 support constraint (3) under the complex behavior policy  $\pi_\beta$  into a simple latent-space constraint:

$$\text{supp}(\pi_L(\cdot|s)) \subseteq \text{supp}(q_L). \quad (6)$$

102 Since  $q_L$  is chosen to have simple support (e.g., a uniform distribution on the  $d$ -ball [7]), this constraint  
 103 is much easier to enforce, thereby decoupling OOD action avoidance from policy optimization.

104 We can thus define the resulting policy optimization problem in the latent space  $\mathcal{L}$  induced by the  
 105 decoder  $\mu_{\theta_{L \rightarrow A}}$ . For a  $Q$  function  $Q_\phi^{\pi_A}$  in the action space (1), this induces a corresponding **latent-space**  
 106  **$Q$  function**  $Q_\phi^{\pi_L}(s, z) : \mathcal{S} \times \mathcal{L} \rightarrow \mathbb{R}$  as the pullback of  $Q_\phi^{\pi_A}$  through the decoder  $\mu_{\theta_{L \rightarrow A}}$ , i.e.,

$$Q_\phi^{\pi_L}(s, z) = Q_\phi^{\pi_A}(s, \mu_{\theta_{L \rightarrow A}}(z; s)). \quad (7)$$

107 As long as the latent-space support constraint (6) is satisfied (e.g., with tanh squashing, clipping or  
 108 projection), *any* off-policy optimization algorithm can then be used to optimize the latent policy  $\pi_L$   
 109 *without* needing to worry about OOD actions, such as DDPG or TD3 in [19], SAC in [8], or even  
 110 plain gradient ascent in [7].

111 **Learning expressive decoders with flow matching.** A key component of latent space RL is the  
 112 decoder  $\mu_{\theta_{L \rightarrow A}}$ . A poorly learned decoder may output actions outside the dataset support and defeat the  
 113 purpose of latent space RL, or be too simple to capture the behavior policy and limit the performance  
 114 of the final policy. Recent works use flow matching [9, 10, 11] to learn  $\mu_{\theta_{L \rightarrow A}}$  due to its ability to  
 115 model complex multimodal behavior distributions [8, 7]. Flow matching learns a time-dependent  
 116 velocity field  $v_{\theta_{L \rightarrow A}}(z, t; s) : \mathcal{L} \times [0, 1] \times \mathcal{S} \rightarrow \mathcal{L}$ , parameterized by  $\theta_{L \rightarrow A}$ , that transports an easy-to-  
 117 sample latent distribution  $q_L \in \Delta(\mathcal{L})$  (e.g., standard Gaussian  $\mathcal{N}(0, I^d)$ ) to the pushforward action  
 118 distribution  $(\mu_{\theta_{L \rightarrow A}}^s)_\# q_L \approx \pi_\beta(\cdot|s)$  over the time interval  $t \in [0, 1]$ , using the flow-matching loss:

$$L_{L \rightarrow A}(\theta_{L \rightarrow A}) = \mathbb{E}_{(s, a) \sim \mathcal{D}, z \sim q_L, t \sim \mathcal{U}[0, 1]} [\|v_{\theta_{L \rightarrow A}}(x_t, t; s) - (a - z)\|^2], \quad x_t = (1 - t)z + ta. \quad (8)$$

119 The corresponding **BC flow map**  $\psi_{\theta_{L \rightarrow A}}(z, t; s) : \mathcal{L} \times [0, 1] \times \mathcal{S} \rightarrow \mathcal{L}$  for the velocity field  $v_{\theta_{L \rightarrow A}}$  can  
 120 be found as the solution to the ODE with initial condition  $z$  at  $t = 0$  and velocity field  $v_{\theta_{L \rightarrow A}}$ , i.e.,

$$\psi_{\theta_{L \rightarrow A}}(z, t; s) = z + \int_0^t v_{\theta_{L \rightarrow A}}(\psi_{\theta_{L \rightarrow A}}(z, \tau; s), \tau; s) d\tau. \quad (9)$$

121 The decoder  $\mu_{\theta_{L \rightarrow A}}$  is then defined as the BC flow map at time  $t = 1$ , i.e.,  $\mu_{\theta_{L \rightarrow A}}(z; s) := \psi_{\theta_{L \rightarrow A}}(z, 1; s)$ .

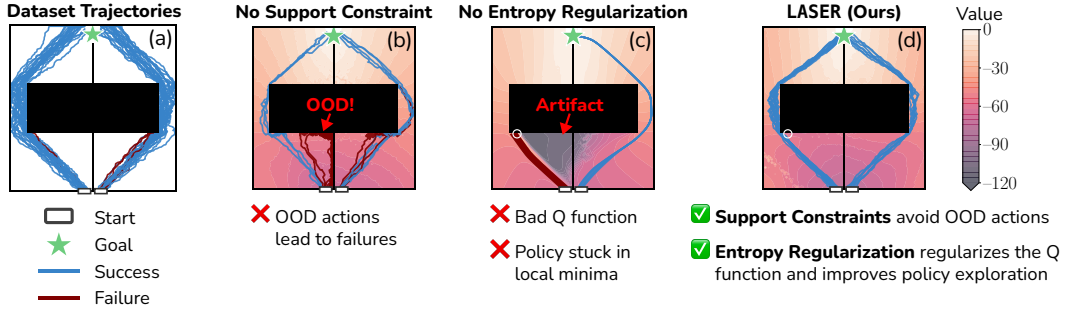


Figure 1: **Both support constraints and entropy regularization are essential.** Without support constraints, the policy generates OOD actions that have unreliable  $Q$ -values and thus fail. Without entropy regularization, the  $Q$ -function exhibits strange artifacts and the policy gets stuck in local minima, resulting in failures on the left side. By combining support constraints and entropy regularization, LASER avoids OOD actions, regularizes the  $Q$ -function, and encourages policy exploration, resulting in a successful policy.

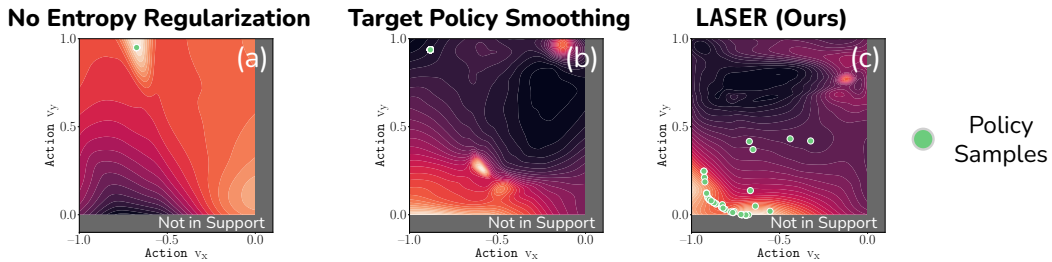


Figure 2:  **$Q$ -function and policy are poorly learned without entropy regularization.** We plot the learned  $Q$ -functions (lighter colors denote higher value) at the state marked by the white circle in Figure 1c. Without entropy regularization, the policy becomes nearly deterministic and clusters at a **false** maximum which causes the policy to move up and get stuck at the wall (Figure 1c). Adding noise to the policy in TD learning (target policy smoothing [20]) fixes the  $Q$ -function landscape, but the policy remains nearly deterministic and is suboptimal. With entropy regularization, the policy samples spread out more and yield a better regularized  $Q$ -function, enabling the policy to move around the wall and reach the target (Figure 1d).

### 122 3 Latent Space Adjoint Matching for Support Constrained Entropy 123 Regularized Offline RL (LASER)

#### 124 3.1 Entropy Regularization is Essential for Latent Space RL

125 Given that latent space RL methods claim to solve the OOD problem, one might be tempted to  
126 conclude that directly performing policy optimization in the latent space with support constraints  
127 results in a strong policy. However, we find this to be false. We illustrate this on a simple 2D  
128 reach-avoid environment (see Appendix E for details). Choosing  $\pi_L$  to be a flow policy as in [7], i.e.,  
129 for a base distribution  $q_B$  and velocity field  $v_{\theta_{B-L}}$ , and corresponding flow map  $\psi_{\theta_{B-L}}$ ,

$$\pi_L(\cdot|s) = (\mu_{\theta_{B-L}}^s)_{\#} q_B, \quad \mu_{\theta_{B-L}}^s(w) := \psi_{\theta_{B-L}}(w, 1), \quad (10)$$

130 we find that naively training  $\pi_L$  to maximize  $Q_{\phi}^{\pi_L}$  with reparametrized gradients and enforcing the  
131 latent-space support constraint (6) via tanh squashing (Figure 1c) gets stuck in local minima and fails  
132 to reach the target (Figure 1). Moreover, the value function is poorly learned and exhibits artifacts.  
133 We explore this further by visualizing the learned  $Q$  function and policy  $\pi_A$  at the state right before  
134 the policy gets stuck at the wall (Figure 2a). Surprisingly, the learned policy is maximizing the  
135 learned  $Q$ -function and is not stuck in a local optimum despite the policy failing to reach the target.  
136 Using target policy smoothing (TPS) [20] to smooth the learned  $Q$ -function by adding noise to the  
137 target action  $a'$  in the TD learning target (2) fixes the  $Q$ -function landscape (Figure 2b), but the  
138 resulting policy is now suboptimal. This is a problem of the policy having insufficient exploration,  
139 which can be mitigated by increasing policy entropy [12].

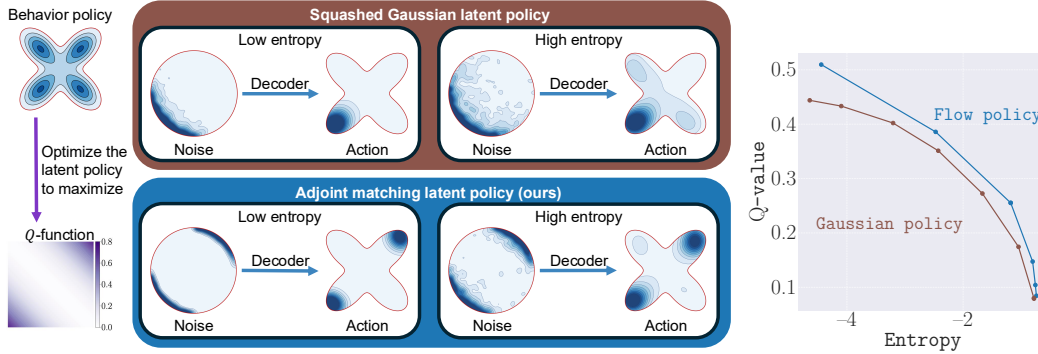


Figure 3: **Flow latent policy + adjoint matching Pareto-dominates a squashed Gaussian latent policy.** For every entropy and  $Q$ -value achieved by the Gaussian latent policy, there exists an operating point achievable by the flow latent policy that has both higher entropy and higher  $Q$ -value.

140 This motivates the use of an **explicit entropy regularization term** to obtain a policy with higher  
 141 entropy, which not only helps with policy optimization via exploration, but also obviates the need  
 142 for TPS in TD learning and yields a better regularized  $Q$ -function. Namely, for weight  $\alpha > 0$ , we  
 143 consider the following entropy-regularized objective to train the latent-space policy  $\pi_L$ :

$$L_{B \rightarrow L}(\theta_{B \rightarrow L}) = -\mathbb{E}_{s \sim \mathcal{D}} [\mathbb{E}_{w \sim q_B} [Q^{\pi_L}(s, \mu_{\theta_{B \rightarrow L}}(w; s))] + \alpha \mathcal{H}(\mu_{\theta_{B \rightarrow L}}(\cdot; s))]. \quad (11)$$

144 The entropy term can also be optimized using reparametrized gradients by computing the log  
 145 probability under the flow map  $\mu_{\theta_{B \rightarrow L}}$  using the continuity equation [9] to obtain

$$\mathcal{H}(\mu_{\theta_{B \rightarrow L}}(\cdot; s)) = \mathbb{E}_{w \sim q_B} \left[ -\log q_B(w) + \int_0^1 \nabla_{x_t} \cdot v_{\theta_{B \rightarrow L}}(x_t, t; s) dt \right], \quad x_t = \psi_{\theta_{B \rightarrow L}}(w, t; s). \quad (12)$$

146 We discuss a more efficient method that does not require backpropagation through the flow map  
 147 in the next subsection. Adding explicit entropy regularization achieves both improved  $Q$ -function  
 148 and improved policy (Figure 2c), resulting in a policy that successfully moves around the wall and  
 149 reaches the target in the simple illustrative environment (Figure 1d).

150 *Remark 3.1.* Examining recent latent space RL methods, we find that their policy extraction does  
 151 include an entropy regularization term despite its importance not being made explicit in their papers.  
 152 This is done either explicitly, as in DSRL [8] which uses a Gaussian latent policy and an SAC-  
 153 style objective [12], or implicitly as in ReFORM [7], where the geometric reflection term prevents  
 154 latent collapse (see Appendix C.1). However, using an expressive flow policy for the latent-space  
 155 policy and explicitly considering entropy regularization enables our method to empirically Pareto-  
 156 dominate Gaussian policies (Figure 3) and maximize the  $Q$ -function better than other implicit entropy  
 157 regularization schemes (Appendix C.1).

158 **Why use flow policy for  $\pi_L$ ?** While a Gaussian latent-space policy significantly simplifies entropy-  
 159 regularized optimization, we motivate our use of a flow model with two arguments. (1) A Gaussian  
 160  $\pi_L$  produces unimodal latents, which restricts the ability for the full policy in action space to capture  
 161 multimodality [7]. (2) While high-entropy policies facilitate robust  $Q$ -function learning and help  
 162 avoid local minima, they often yield lower  $Q$ -values. By being more expressive, flow policies  
 163 empirically Pareto-dominate Gaussian policies on the performance-entropy tradeoff (Figure 3).

### 164 3.2 Efficient Entropy-Regularized Policy Optimization via Adjoint Matching

165 While (11) can be used as-is to train the latent policy  $\pi_L$  with reparametrized gradients, computing the  
 166 gradient requires backpropagating through the flow map  $\mu_{\theta_{B \rightarrow L}}$  (9) and the entropy (12), both of which  
 167 are solutions to ODEs and require backpropagation through time (BPTT), which is computationally  
 168 expensive and prone to numerical instability [21, 22, 23].

169 To circumvent these computational bottlenecks, we adopt an alternative approach. Instead of optimiz-  
 170 ing the loss directly via BPTT, we derive the optimal latent-space policy  $\pi_L^*(\cdot|s) = (\mu_{\theta_{B \rightarrow L}}^s)^* q_B$  that  
 171 minimizes the entropy-regularized policy loss (11) in the typical case where the support of the latent  
 172 space prior distribution  $q_L$  has finite Lebesgue measure (e.g.,  $d$ -ball).

173 **Lemma 3.2.** Fix a state  $s \in \mathcal{S}$ . Suppose  $\text{supp}(q_L(\cdot|s))$  has finite Lebesgue measure, and the  $Q$ -  
 174 function  $Q_\phi^{\pi_L}(s, \cdot)$  has a finite upper bound on  $\text{supp}(q_L(\cdot|s))$ . Then, the optimal distribution  $\pi_L^*$  that  
 175 minimizes the loss (11) is given by the following  $Q$ -tilted distribution:

$$\pi_L^*(\cdot|s) \propto q_L^{\text{unif}}(\cdot) \exp\left(Q_\phi^{\pi_L}(s, \cdot)/\alpha\right), \quad (13)$$

176 where  $q_L^{\text{unif}} = \mathcal{U}(\text{supp}(q_L(\cdot|s)))$  is the uniform distribution on the support of  $q_L$ .

177 See Appendix B for the proof. With Lemma 3.2, we can optimize (11) by learning to sample from the  
 178 optimal distribution  $\pi_L^*$ , which can be done efficiently *without BPTT* using **adjoint matching** [16].

179 Namely, suppose both  $\pi_L$  and  $q_L^{\text{unif}}$  are flow models with a common base distribution  $q_B$ , i.e., (10)  
 180 holds for  $\pi_L$ , and for flow map  $\psi_{B \rightarrow L}^{\text{base}}$  generated by velocity field  $v_{B \rightarrow L}^{\text{base}}$ ,

$$q_L^{\text{unif}} = (\mu_{B \rightarrow L}^{\text{base}})_\# q_B, \quad \mu_{B \rightarrow L}^{\text{base}}(w) = \psi_{B \rightarrow L}^{\text{base}}(w, 1), \quad (14)$$

181 and let  $p_v(\{w_t\}_{t \in [0,1]}) \in \Delta(\mathcal{B}^{[0,1]})$  denote the distribution over the flow trajectory  $\{w_t\}_{t \in [0,1]}$  under  
 182 velocity  $v$ . Moreover, consider the tilted distribution over *trajectories*  $p_{v^*}$ , defined as

$$p_{v^*}(\{w_t\}_{t \in [0,1]}) \propto p_{v_{B \rightarrow L}^{\text{base}}}(\{w_t\}_{t \in [0,1]}) \exp\left(Q_\phi^{\pi_L}(s, w_1)/\alpha\right), \quad (15)$$

183 where (13) now corresponds to the marginal of  $p_{v^*}$  at  $t = 1$ . Then,  $p_{v^*}$  can be learned without BPTT  
 184 by repeatedly optimizing a regression loss using a “special estimate” of  $v^*$  as the regression target.  
 185 Specifically, the target distribution over trajectories (15) corresponding to the vector field  $v^*$  is *also*  
 186 the distribution of the following stochastic differential equation (SDE)

$$dw_t = (2v(w_t, t; s) - w_t/t) dt + \sigma_t dB_t, \quad (16)$$

187 for a memoryless noise schedule  $\sigma_t$ , standard Brownian motion  $B_t$ , and a control  $v$  that solves the  
 188 following stochastic optimal control (SOC) problem with quadratic control cost and terminal cost  
 189 given by the negative  $Q$ -function [16]:

$$\min_v \mathbb{E}_{s \sim \mathcal{D}, \{w_t\}} \left[ \int_0^1 \frac{2}{\sigma_t^2} \|v(w_t, t; s) - v_{B \rightarrow L}^{\text{base}}(w_t, t)\|^2 dt - Q_\phi^{\pi_L}(s, w_1)/\alpha \right]. \quad (17)$$

190 Then, given a rollout  $\{w_t\}_{t \in [0,1]}$  sampled with control  $v$ , we can construct the following estimate  $\hat{v}^*$   
 191 of the optimal control  $v^*$  as

$$\hat{v}^*(w_t, t; s) = v_{B \rightarrow L}^{\text{base}}(w_t, t) - \sigma_t^2 \tilde{g}_t / 2, \quad (18)$$

192 where  $\tilde{g}_t$  is the “lean” adjoint state computed by integrating the reverse ODE using  $\{w_t\}_{t \in [0,1]}$

$$d\tilde{g}_t = -\nabla_{w_t} [2v_{B \rightarrow L}^{\text{base}}(w_t, t) - w_t/t] \tilde{g}_t dt, \quad \tilde{g}_1 = -\nabla_{w_1} Q_\phi^{\pi_L}(s, w_1)/\alpha. \quad (19)$$

193 Using the estimate  $\hat{v}^*$  as the regression target in the following regression loss then results in a unique  
 194 fixed point at the optimal control  $v^*$  [16], where we replace the current velocity field  $v$  with the  
 195 trainable velocity field  $v_{\theta_{B \rightarrow L}}$  to be optimized:

$$L_{\text{AM}}(\theta_{B \rightarrow L}) = \mathbb{E}_{s \sim \mathcal{D}, \{w_t\}} \left[ \frac{1}{2} \int_0^1 \left\| \frac{2}{\sigma_t} \left( v_{\theta_{B \rightarrow L}}(w_t, t; s) - \hat{v}^*(w_t, t; s) \right) \right\|^2 dt \right]. \quad (20)$$

196 Importantly, optimizing (20) **does not require backpropagation through the flow** and only requires  
 197 the solution of an ODE to compute the adjoint state  $\tilde{g}_t$  for the regression target  $\hat{v}^*$ , circumventing  
 198 the computational and numerical stability issues of BPTT. For a more comprehensive mathematical  
 199 treatment of adjoint matching, we refer readers to [16].

### 200 3.3 Algorithm

201 We now present **Latent Space Adjoint Matching for Support Constrained Entropy Regularized**  
 202 **Offline RL (LASER)**. LASER combines the aforementioned insights to perform latent space RL with an  
 203 expressive flow decoder and efficient entropy-regularized policy optimization via adjoint matching.

- 204 • Following [7], we choose the latent prior  $q_L = \mathcal{U}(\mathcal{B}_l^d)$  to be the uniform distribution on the ball  $\mathcal{B}_l^d$   
 205 of radius  $l$  in  $\mathbb{R}^d$ , which has finite Lebesgue measure and satisfies the conditions of Lemma 3.2.

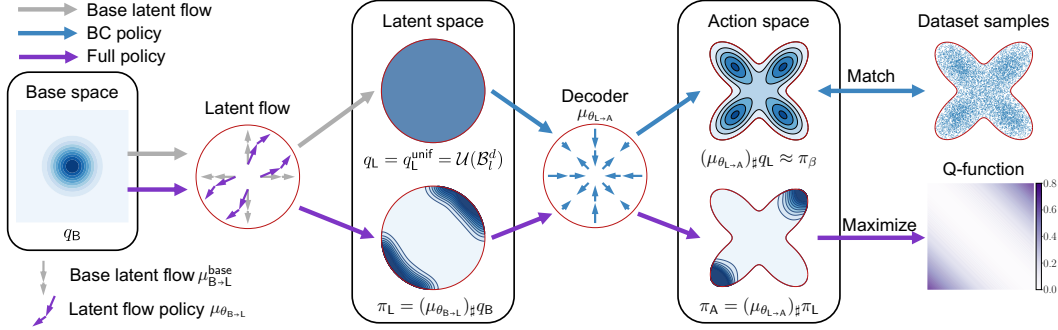


Figure 4: **LASER Algorithm.** ( $\rightarrow$ ) We train a flow-matching decoder  $\mu_{\theta_{L \rightarrow A}}$  that maps samples from the uniform distribution in the  $d$ -ball  $q_L = q_L^{\text{unif}} = \mathcal{U}(\mathcal{B}_L^d)$  in *latent space* to the dataset actions in *action space*. ( $\rightarrow$ ) We then learn a base latent flow  $\mu_{B \rightarrow L}^{\text{base}}$  that maps the prior distribution  $q_B$  in the base space to the uniform distribution in the latent space  $q_L^{\text{unif}}$ . ( $\rightarrow$ ) The base latent flow is then fine-tuned via adjoint matching, which results in a latent flow policy  $\mu_{\theta_{B \rightarrow L}}$  that generates the final action distribution  $\pi_A = (\mu_{\theta_{L \rightarrow A}})_{\#} (\mu_{\theta_{B \rightarrow L}})_{\#} q_B$ .

206 • Adjoint matching requires  $q_L^{\text{unif}}$ , the uniform distribution over the support of  $q_L$ , to also be a flow  
 207 model. While an analytical vector field can be derived since  $q_L^{\text{unif}} = q_L$  is the uniform distribution  
 208 over the ball, it involves special functions with no easy closed-form expression (see Appendix C.2).  
 209 Instead, we train a flow model via flow-matching to match  $q_L^{\text{unif}}$  with the same base distribution  $q_B$   
 210 as  $\pi_L$  to simplify the implementation.

211 We now summarize LASER (Figure 4), which consists of four components that happen simultaneously:  
 212 **(1. Flow-matching decoder)** We train a state-conditioned flow decoder  $\mu_{\theta_{L \rightarrow A}}$  using flow-matching  
 213 to map samples from the uniform distribution on the  $d$ -ball  $q_L$  to dataset actions, i.e., behavior  
 214 cloning (Section 2.2). **(2. Flow-matching for  $q_L^{\text{unif}}$ )** We train a flow model  $\mu_{B \rightarrow L}^{\text{base}}$  to match the  
 215 uniform distribution  $q_L^{\text{unif}}$  in latent space from  $q_B$  to use in adjoint matching. **(3. Critic learning)**  
 216 We learn a critic  $Q_{\phi}^{\mu_{\theta}}$  using the standard TD learning objective (2). **(4. Entropy-regularized  
 217 policy optimization)** Finally, we perform entropy-regularized policy optimization in latent space  
 218 to maximize the latent-space  $Q$ -function  $Q_{\phi}^{\pi_L}$  (7) using adjoint matching (Section 3.2). We provide  
 219 practical implementation details and a high-level algorithm for LASER in Appendix D.

## 220 4 Related Work

221 We discuss the most closely related work in this section and defer a broader discussion of offline RL  
 222 and generative policies to Appendix A.

223 **Distribution shift in offline RL.** A central challenge in offline RL is distribution shift: learned value  
 224 functions can overestimate actions outside the dataset support, causing policy improvement to select  
 225 out-of-distribution actions with poor true value. Prior methods address this issue by regularizing  
 226 the policy toward the behavior policy with statistical distances, learning pessimistic value functions,  
 227 or explicitly restricting policy improvement to supported actions [5, 24, 25, 26, 4]. We follow the  
 228 support-constrained perspective, using an expressive flow policy to enforce support by construction.

229 **Expressive generative policies for offline RL.** Diffusion [27, 28] and flow-matching [9, 11, 10]  
 230 models have recently been used as expressive alternatives to deterministic or Gaussian policies,  
 231 especially in robotic imitation learning where action distributions are often multimodal [29, 30].  
 232 These policies have also been adapted to offline RL, where the goal is to improve beyond the  
 233 behavior policy using a learned critic. Existing methods differ in how they use the critic: some  
 234 use critic values to select, reweight, or clone high-value actions [31, 32, 33, 34], while others use  
 235 critic gradients to directly improve the generated actions [35, 30, 21]. Most closely related to our  
 236 optimization mechanism, [23] fine-tunes flow policies using critic gradients under a KL-regularized  
 237 reward objective. Unlike these methods, which rely on statistical regularization to control distribution  
 238 shift, LASER enforces support constraints through its latent action parameterization.

239 **Support-constrained policy extraction with latent-space RL.** Recent methods avoid OOD actions  
240 by learning a latent action space whose decoder maps into the support of the behavior policy, then  
241 performing policy optimization in this latent space [19, 8, 7]. Our method builds on this idea, but  
242 revisits the optimization problem that remains once support has been controlled. In particular, we  
243 show that entropy regularization is important for preventing policy collapse and stabilizing policy  
244 extraction within the supported action set.

245 **Entropy regularization.** Entropy regularization is a classical tool for smoothing policy optimization,  
246 improving robustness, and encouraging exploration [13, 12, 15]. In offline RL, though some works  
247 have noted its benefits for exploration during online fine-tuning [35, 36], it receives much less  
248 attention, perhaps due to the implicit entropy regularization effect from behavior policy regularization  
249 techniques. Our perspective is different: because support constraints are enforced by the policy  
250 parameterization, entropy regularization can be used for its standalone role in policy extraction.

## 251 5 Experiments

252 We conduct experiments to answer the following research questions. **(Q1):** How does LASER  
253 perform against other offline RL algorithms with flow policies? **(Q2):** How sensitive is LASER to  
254 hyperparameters? **(Q3):** How does adjoint matching compare with BPTT? **(Q4):** Does the latent  
255 space decouple entropy regularization from behavior regularization? Details on implementation,  
256 baselines, environments, and additional results and ablations are in Appendix F.

### 257 5.1 Setup

258 **Environments.** We evaluate our method on the OGBench offline RL benchmark [18]. The evalua-  
259 tion suite spans four environments, covering both locomotion and manipulation domains, with five  
260 tasks per environment. Each task is evaluated on two distinct datasets (CLEAN and NOISY), yielding  
261 a total of  $4 \times 5 \times 2 = 40$  evaluation settings.

262 **Baselines.** Since recent works consistently demonstrate flow-based policies outperform traditional  
263 alternatives [21, 37, 22, 23], we benchmark LASER against state-of-the-art flow-based methods. These  
264 include IFQL [32, 21], which fine-tunes a BC flow policy via  $Q$ -value-weighted importance sampling;  
265 FQL [21], which applies Wasserstein distance regularization; and QAM [23], which uses adjoint  
266 matching to approximate the  $Q$ -value-weighted BC action distribution, alongside its extension QAM-E  
267 [23] that learns an additional residual policy. We also compare against latent space RL methods,  
268 specifically DSRL [8] with a Gaussian latent policy, and ReFORM [7] with a reflected flow latent policy.  
269 Crucially, unlike the other baselines using fine-tuned hyperparameters per environment and dataset,  
270 LASER, ReFORM, and IFQL use a *constant set of hyperparameters across all evaluated tasks*.

### 271 5.2 Results

272 **(Q1): LASER achieves state-of-the-art performance using a single set of hyperparameters.** In  
273 Figure 5, we present the performance profiles [38] aggregating all tasks across both the CLEAN and  
274 NOISY datasets. The profiles demonstrate that LASER dominates all evaluated baselines. Notably,  
275 the two support constraint methods (LASER and ReFORM) consistently outperform other approaches  
276 with the NOISY dataset, underscoring the advantage of enforcing support constraints, which does not  
277 restrict policy improvement. Furthermore, LASER significantly improves upon ReFORM. By avoiding  
278 the gradient vanishing issues inherent to the design of ReFORM’s reflection term, LASER is able to  
279 achieve a perfect success rate of 1.0 on nearly half of the evaluated tasks.

280 **(Q2): The key hyperparameter of LASER is the inverse temperature.** Using a constant set of  
281 hyperparameters for LASER across all tasks demonstrates its robustness to different tasks and datasets.  
282 Nevertheless, we analyze the sensitivity of LASER to the inverse temperature  $1/\alpha$  (Figure 6). Values of  
283  $1/\alpha \in [5, 20]$  yield comparable results, but performance drops outside this range from too much/too  
284 little entropy. Nevertheless, we find that  $1/\alpha = 10$  yields stable results for all environments. We also  
285 analyze sensitivity to  $\sigma_0$  and find that LASER is robust to this hyperparameter (Appendix F.5).

286 **(Q3): Adjoint matching achieves better results and is faster than BPTT using reparametrized**  
287 **gradients.** A key advantage of adjoint matching is that it circumvents the need to perform BPTT

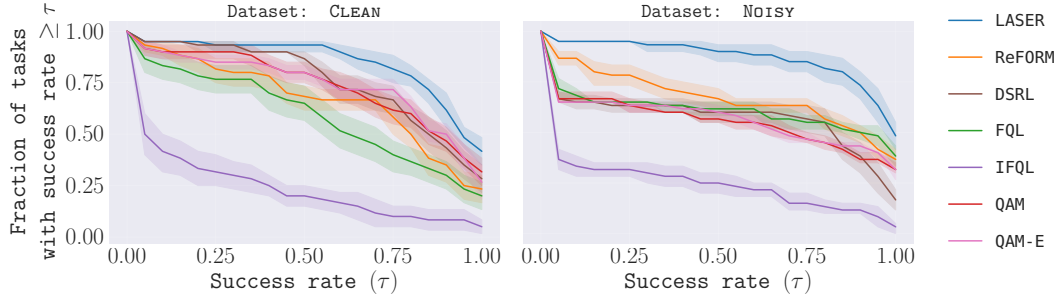


Figure 5: **Performance profile plots for all tasks with different datasets.** For a given success rate  $\tau$  (x-axis), we plot the fraction of tasks with a success rate  $\geq \tau$  [38]. For both datasets, **LASER** achieves higher success rates across a larger fraction of tasks than all other baselines.

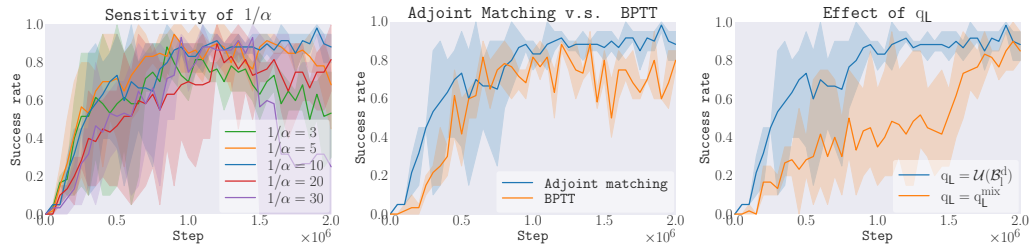


Figure 6: **Ablation studies in cube-double-play-singletask-task2-v0.** Moderate values of the inverse temperature  $1/\alpha \in [5, 20]$  yield comparable results. Adjoint matching achieves better and more stable results than BPTT. Using a different  $q_L$  does not affect the final performance of **LASER**, indicating that the latent space successfully decouples entropy regularization from behavior regularization.

288 on (12), which is computationally expensive and numerically unstable. We validate this empirically  
 289 and observe that adjoint matching yields more stable training and stronger performance compared to  
 290 BPTT (Figure 6) while training  $3.6\times$  faster (123 it/s vs 34 it/s).

291 **(Q4): The latent space successfully decouples entropy regularization from behavior regular-**  
 292 **ization.** To understand the improvements of **LASER**, we change  $q_L$  to a distribution with the same  
 293 support but a different density, which changes the behavior regularization effect but keeps the same  
 294 latent space entropy regularization. While convergence is slower, the final performance is unchanged,  
 295 suggesting that the improvements of **LASER** come from entropy regularization (Figure 6).

## 296 6 Conclusion

297 We introduced **LASER**, an entropy-regularized approach to support-constrained flow policy opti-  
 298 mization for offline RL. Our key insight is that support constraints and entropy regularization solve  
 299 different problems: the former prevents OOD actions, while the latter stabilizes policy extraction  
 300 within the learned dataset support. By combining a flow-based decoder with latent space adjoint  
 301 matching, **LASER** effectively prevents mode collapse, smooths the optimization landscape, and avoids  
 302 local minima, demonstrating state-of-the-art performance on the OGBench suite using a fixed set of  
 303 hyperparameters, eliminating the need for per-task tuning common in prior offline RL methods.

304 **Limitations and Future Work.** We identify several promising avenues for extension. First,  
 305 although adjoint matching circumvents the need for BPTT, solving the ODEs needed to sample from  
 306 the flow model during training remains computationally intensive for deep neural ODEs, which could  
 307 be accelerated with distillation techniques [22]. Second, by relying on the decoder to enforce support  
 308 constraints, **LASER** inherits OOD errors made by the decoder itself. Integrating behavior cloning  
 309 methods with stricter support constraints, diagnosing when the BC model generates OOD errors, or  
 310 applying a pre-trained robust BC model could mitigate this. In addition, **LASER** applies the simplest  
 311 value function learning method and actor-critic structure similar to [21], which could be improved  
 312 upon [39, 40, 41, 37].

## References

- [1] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- [2] Qing Wang, Jiechao Xiong, Lei Han, Han Liu, Tong Zhang, et al. Exponentially weighted imitation learning for batched historical data. *Advances in Neural Information Processing Systems*, 31, 2018.
- [3] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- [4] Yixiu Mao, Hongchang Zhang, Chen Chen, Yi Xu, and Xiangyang Ji. Supported trust region optimization for offline reinforcement learning. In *International Conference on Machine Learning*, pages 23829–23851. PMLR, 2023.
- [5] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. *Advances in neural information processing systems*, 32, 2019.
- [6] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- [7] Songyuan Zhang, Oswin So, H M Sabbir Ahmad, Eric Yang Yu, Matthew Cleaveland, Mitchell Black, and Chuchu Fan. ReFORM: Reflected flows for on-support offline RL via noise manipulation. In *The Fourteenth International Conference on Learning Representations*, 2026.
- [8] Andrew Wagenmaker, Mitsuhiro Nakamoto, Yunchu Zhang, Seohong Park, Waleed Yagoub, Anusha Nagabandi, Abhishek Gupta, and Sergey Levine. Steering your diffusion policy with latent space reinforcement learning. *arXiv preprint arXiv:2506.15799*, 2025.
- [9] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- [10] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023.
- [11] Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *arXiv preprint arXiv:2209.15571*, 2022.
- [12] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. Pmlr, 2018.
- [13] Brian D Ziebart. *Modeling purposeful adaptive behavior with the principle of maximum causal entropy*. Carnegie Mellon University, 2010.
- [14] Benjamin Eysenbach and Sergey Levine. Maximum entropy RL (provably) solves some robust RL problems. In *International Conference on Learning Representations*, 2022.
- [15] Zafarali Ahmed, Nicolas Le Roux, Mohammad Norouzi, and Dale Schuurmans. Understanding the impact of entropy on policy optimization. In *International conference on machine learning*, pages 151–160. PMLR, 2019.
- [16] Carles Domingo-Enrich, Michal Drozdal, Brian Karrer, and Ricky TQ Chen. Adjoint matching: Fine-tuning flow and diffusion generative models with memoryless stochastic optimal control. *arXiv preprint arXiv:2409.08861*, 2024.
- [17] Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992.

- 358 [18] Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking  
359 offline goal-conditioned rl. In *International Conference on Learning Representations (ICLR)*,  
360 2025.
- 361 [19] Wenxuan Zhou, Sujay Bajracharya, and David Held. Plas: Latent action space for offline  
362 reinforcement learning. In *Conference on Robot Learning*, pages 1719–1735. PMLR, 2021.
- 363 [20] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error  
364 in actor-critic methods. In *International conference on machine learning*, pages 1587–1596.  
365 PMLR, 2018.
- 366 [21] Seohong Park, Qiyang Li, and Sergey Levine. Flow q-learning. In *International Conference on*  
367 *Machine Learning (ICML)*, 2025.
- 368 [22] Nicolas Espinosa-Dice, Yiyi Zhang, Yiding Chen, Bradley Guo, Owen Oertell, Gokul Swamy,  
369 Kianté Brantley, and Wen Sun. Scaling offline RL via efficient and expressive shortcut models.  
370 In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- 371 [23] Qiyang Li and Sergey Levine. Q-learning with adjoint matching. In *The Fourteenth International*  
372 *Conference on Learning Representations*, 2026.
- 373 [24] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for of-  
374 fline reinforcement learning. *Advances in neural information processing systems*, 33:1179–1191,  
375 2020.
- 376 [25] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit  
377 q-learning. In *International Conference on Learning Representations*, 2022.
- 378 [26] Jialong Wu, Haixu Wu, Zihan Qiu, Jianmin Wang, and Mingsheng Long. Supported policy  
379 optimization for offline reinforcement learning. *Advances in Neural Information Processing*  
380 *Systems*, 35:31278–31291, 2022.
- 381 [27] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and  
382 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*  
383 *preprint arXiv:2011.13456*, 2020.
- 384 [28] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*  
385 *in neural information processing systems*, 33:6840–6851, 2020.
- 386 [29] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ  
387 Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion.  
388 *The International Journal of Robotics Research*, page 02783649241273668, 2023.
- 389 [30] Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive  
390 policy class for offline reinforcement learning. In *The Eleventh International Conference on*  
391 *Learning Representations*, 2023.
- 392 [31] Huayu Chen, Cheng Lu, Chengyang Ying, Hang Su, and Jun Zhu. Offline reinforcement learning  
393 via high-fidelity generative behavior modeling. In *The Eleventh International Conference on*  
394 *Learning Representations*, 2023.
- 395 [32] Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey  
396 Levine. Idql: Implicit q-learning as an actor-critic method with diffusion policies. *arXiv preprint*  
397 *arXiv:2304.10573*, 2023.
- 398 [33] Bingyi Kang, Xiao Ma, Chao Du, Tianyu Pang, and Shuicheng Yan. Efficient diffusion  
399 policies for offline reinforcement learning. *Advances in Neural Information Processing Systems*,  
400 36:67195–67212, 2023.
- 401 [34] Shiyuan Zhang, Weitong Zhang, and Quanquan Gu. Energy-weighted flow matching for  
402 offline reinforcement learning. In *Proceedings of the International Conference on Learning*  
403 *Representations (ICLR 2025)*, 2025.

- 404 [35] Shutong Ding, Ke Hu, Zhenhao Zhang, Kan Ren, Weinan Zhang, Jingyi Yu, Jingya Wang, and  
405 Ye Shi. Diffusion-based reinforcement learning via q-weighted variational policy optimization.  
406 *Advances in Neural Information Processing Systems*, 37:53945–53968, 2024.
- 407 [36] Perry Dong, Qiyang Li, Dorsa Sadigh, and Chelsea Finn. Expo: Stable reinforcement learning  
408 with expressive policies. *arXiv preprint arXiv:2507.07986*, 2025.
- 409 [37] Bhavya Agrawalla, Michal Nauman, Khush Agrawal, and Aviral Kumar. floq: Training critics  
410 via flow-matching for scaling compute in value-based rl. *arXiv preprint arXiv:2509.06863*,  
411 2025.
- 412 [38] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Belle-  
413 mare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural  
414 information processing systems*, 34:29304–29320, 2021.
- 415 [39] Yixiu Mao, Hongchang Zhang, Chen Chen, Yi Xu, and Xiangyang Ji. Supported value  
416 regularization for offline reinforcement learning. *Advances in Neural Information Processing  
417 Systems*, 36:40587–40609, 2023.
- 418 [40] Divyansh Garg, Joey Hejna, Matthieu Geist, and Stefano Ermon. Extreme q-learning: Maxent  
419 rl without entropy. In *International Conference on Learning Representations*, 2023.
- 420 [41] Tenglong Liu, Yang Li, Yixing Lan, Hao Gao, Wei Pan, and Xin Xu. Adaptive advantage-guided  
421 policy regularization for offline reinforcement learning. In *International Conference on Machine  
422 Learning*, pages 31406–31424. PMLR, 2024.
- 423 [42] Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza,  
424 Noah Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch deep reinforcement  
425 learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- 426 [43] Noah Y Siegel, Jost Tobias Springenberg, Felix Berkenkamp, Abbas Abdolmaleki, Michael  
427 Neunert, Thomas Lampe, Roland Hafner, Nicolas Heess, and Martin Riedmiller. Keep doing  
428 what worked: Behavioral modelling priors for offline reinforcement learning. *arXiv preprint  
429 arXiv:2002.08396*, 2020.
- 430 [44] Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. Awac: Accelerating online  
431 reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- 432 [45] Ziyu Wang, Alexander Novikov, Konrad Zolna, Josh S Merel, Jost Tobias Springenberg, Scott E  
433 Reed, Bobak Shahriari, Noah Siegel, Caglar Gulcehre, Nicolas Heess, et al. Critic regularized  
434 regression. *Advances in Neural Information Processing Systems*, 33:7768–7778, 2020.
- 435 [46] Jiafei Lyu, Xiaoteng Ma, Xiu Li, and Zongqing Lu. Mildly conservative q-learning for offline  
436 reinforcement learning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun  
437 Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- 438 [47] Zhepeng Cen, Zuxin Liu, Zitong Wang, Yihang Yao, Henry Lam, and Ding Zhao. Learning  
439 from sparse offline datasets via conservative density estimation. In *The Twelfth International  
440 Conference on Learning Representations*, 2024.
- 441 [48] Qingjun Wang, Hongtu Zhou, Hang Yu, Junqiao Zhao, Yanping Zhao, Chen Ye, Ziqiao Wang,  
442 and Guang Chen. Beyond penalization: Diffusion-based out-of-distribution detection and  
443 selective regularization in offline reinforcement learning. In *The Fourteenth International  
444 Conference on Learning Representations*, 2026.
- 445 [49] Jing Zhang, Chi Zhang, Wenjia Wang, and Bingyi Jing. Constrained policy optimization with  
446 explicit behavior density for offline reinforcement learning. *Advances in Neural Information  
447 Processing Systems*, 36:5616–5630, 2023.
- 448 [50] Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal conditioned imitation  
449 learning using score-based diffusion policies. In *Robotics: Science and Systems*, 2023.
- 450 [51] Tim Pearce, Tabish Rashid, Anssi Kanervisto, Dave Bignell, Mingfei Sun, Raluca Georgescu,  
451 Sergio Valcarcel Macua, Shan Zheng Tan, Ida Momennejad, Katja Hofmann, et al. Imitating  
452 human behaviour with diffusion models. *arXiv preprint arXiv:2301.10677*, 2023.

- 453 [52] Huayu Chen, Cheng Lu, Zhengyi Wang, Hang Su, and Jun Zhu. Score regularized policy  
454 optimization through diffusion behavior. In *The Twelfth International Conference on Learning*  
455 *Representations*, 2024.
- 456 [53] Zihan Ding and Chi Jin. Consistency models as a rich and efficient policy class for reinforcement  
457 learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- 458 [54] Ruoqi Zhang, Ziwei Luo, Jens Sjölund, Thomas B Schön, and Per Mattsson. Entropy-regularized  
459 diffusion policy with q-ensembles for offline reinforcement learning. *Advances in neural*  
460 *information processing systems*, 37:98871–98897, 2024.
- 461 [55] Suzan Ece Ada, Erhan Oztop, and Emre Ugur. Diffusion policies for out-of-distribution  
462 generalization in offline reinforcement learning. *IEEE Robotics and Automation Letters*,  
463 9(4):3116–3123, 2024.
- 464 [56] Franki Nguimatsia Tiofack, Théotime Le Hellard, Fabian Schramm, Nicolas Perrin-Gilbert, and  
465 Justin Carpentier. Guided flow policy: Learning from high-value actions in offline reinforcement  
466 learning. In *The Fourteenth International Conference on Learning Representations*, 2026.
- 467 [57] Jongseong Chae, Jongeui Park, Yongjae Shin, Gyeongmin Kim, Seungyul Han, and Youngchul  
468 Sung. Flow actor-critic for offline reinforcement learning. In *The Fourteenth International*  
469 *Conference on Learning Representations*, 2026.
- 470 [58] Xiaoyuan Cheng, Haoyu Wang, Wenxuan Yuan, Ziyang Wang, Zonghao Chen, Li Zeng, and  
471 Zhuo Sun. Fisher decorator: Refining flow policy via a local transport map. *arXiv preprint*  
472 *arXiv:2604.17919*, 2026.
- 473 [59] He Zhang, Ying Sun, and Hui Xiong. Goldenstart: Q-guided priors and entropy control for dis-  
474 stilling flow policies. In *The Fourteenth International Conference on Learning Representations*,  
475 2026.
- 476 [60] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis.  
477 *Advances in neural information processing systems*, 34:8780–8794, 2021.
- 478 [61] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*  
479 *arXiv:2207.12598*, 2022.
- 480 [62] Cheng Lu, Huayu Chen, Jianfei Chen, Hang Su, Chongxuan Li, and Jun Zhu. Contrastive  
481 energy prediction for exact energy-guided diffusion sampling in offline reinforcement learning.  
482 In *International Conference on Machine Learning*, pages 22825–22855. PMLR, 2023.
- 483 [63] Kevin Frans, Seohong Park, Pieter Abbeel, and Sergey Levine. Diffusion guidance is a  
484 controllable policy improvement operator. *arXiv preprint arXiv:2505.23458*, 2025.
- 485 [64] Haoran Xu, Kaiwen Hu, Somayeh Sojoudi, and Amy Zhang. Value gradient flow: Behavior-  
486 regularized RL without regularization. In *The Fourteenth International Conference on Learning*  
487 *Representations*, 2026.
- 488 [65] Yilun Du, Conor Durkan, Robin Strudel, Joshua B Tenenbaum, Sander Dieleman, Rob Fergus,  
489 Jascha Sohl-Dickstein, Arnaud Doucet, and Will Sussman Grathwohl. Reduce, reuse, recycle:  
490 Compositional generation with energy-based diffusion models and mcmc. In *International*  
491 *conference on machine learning*, pages 8489–8510. PMLR, 2023.
- 492 [66] Kamyar Ghasemipour, Shixiang Shane Gu, and Ofir Nachum. Why so pessimistic? estimating  
493 uncertainties for offline rl through ensembles, and why their independence matters. *Advances in*  
494 *Neural Information Processing Systems*, 35:18267–18281, 2022.
- 495 [67] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua  
496 Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations*,  
497 *ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- 498 [68] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint*  
499 *arXiv:1606.08415*, 2016.
- 500 [69] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint*  
501 *arXiv:1607.06450*, 2016.

## 502 A Extended Related Work

503 **Mitigating distribution shift in Offline RL.** A central challenge in offline RL is mitigating  
504 distribution shift: learned value functions can overestimate actions that are not in the dataset, causing  
505 the policy to select these out-of-distribution actions. Classical approaches address this issue by either  
506 regularizing the policy to remain close to the behavior policy with statistical distances (e.g., KL  
507 divergence [2, 3, 6, 42, 43, 44, 45, 25], MMD distance [5] or Wasserstein distance [6]), or regularizing  
508 the value function to be pessimistic on out-of-distribution actions [24, 25, 46, 39, 47, 48]. Another  
509 family of methods address the OOD problem by penalizing actions that fall outside the support of  
510 the behavior policy [5, 26, 4, 49]. However, these methods rely on heuristic sample-based proxies  
511 [5] or Gaussian density approximations [26, 4, 49] which can fail to capture expressive multi-modal  
512 behavior policies.

513 **Expressive generative policies for offline RL.** Diffusion [27, 28] and flow-matching [9, 11, 10]  
514 models have recently been used as expressive alternatives to deterministic or Gaussian policies,  
515 especially in robotic imitation learning where action distributions are often multimodal [29, 50, 51, 30].  
516 These policies have also been adapted to offline RL, where the goal is to improve beyond the behavior  
517 policy using a learned critic. However, because flow and diffusion policies generate actions through an  
518 iterative sampling process, critic-based optimization is more challenging than for standard Gaussian  
519 actors. Existing reward-optimization methods for generative policies can be broadly divided by how  
520 they use the learned critic: *critic-evaluation methods*, which use only scalar values  $Q(s, a)$  to select  
521 or reweight actions, and *critic-gradient methods*, which use  $\nabla_a Q(s, a)$  to directly improve generated  
522 actions.

523 (i) *Critic-evaluation* methods include sampling-based methods that imitate the behavior policy and  
524 use the critic to sample from the reward tilted distribution either during test-time [31, 32], or  
525 during training-time as samples to imitate from [52]. These methods often reliably improve the  
526 quality of the extracted policy at the cost of increased computational cost and depend on the  
527 quality of the behavior policy. Weighted behavioral cloning methods also fall under this category.  
528 These methods use the critic to assign per-sample weights in the diffusion or flow-matching loss,  
529 selectively cloning high-value actions in the dataset [33, 35, 34], but generally do not optimize  
530 the learned policy towards the reward-tilted distribution [23].

531 (ii) *Critic-gradient* methods include reparametrized policy gradient methods that have a behavior  
532 cloning term to match the behavior policy and a term that optimizes the critic, which requires  
533 backpropagating through the iterative sampling process [30, 53, 54, 55] and can be computation-  
534 ally heavy and unstable [21, 22, 23], or instead learns a distilled policy [21, 56, 57, 58, 59] at  
535 the expense of policy expressiveness. Another line of work uses guidance [60, 61] during the  
536 iterative sampling process to sample [62, 63, 64], but these methods rely on approximations that  
537 do not have theoretical guarantees [65, 23]. Most closely related to our optimization mecha-  
538 nism, [23] provides an efficient way to fine-tune flow-based policies using critic gradients and a  
539 KL-regularized reward objective via adjoint matching [16]. However, unlike critic-evaluation  
540 methods where samples from the learned policy are either directly cloned (sampling-based  
541 methods) or selectively cloned (weighted BC methods), these methods rely on statistical distance  
542 regularization to avoid the OOD problem (e.g., L2 / Wasserstein distance [21] or KL regulariza-  
543 tion [23]). A key challenge is that the regularization weight is a hyperparameter that significantly  
544 affects the performance and needs to be tuned for each task and dataset ([18, 21]).

545 Our method falls under the *critic-gradient* category as we also use adjoint matching to optimize  
546 the policy with critic gradients under KL regularization. However, to address the OOD problem,  
547 LASER relies on the latent space to enforce the support constraints by construction instead of using  
548 KL regularization as in [23] which can be too restrictive in cases where optimal actions happen with  
549 very low probability under the behavior policy [5].

550 **Support constraints in offline RL.** Instead of addressing the OOD problem with regularization,  
551 another line of work propose to address the OOD problem by learning a bounded latent space that  
552 maps to the support of the behavior policy, then performing policy extraction in this latent space to  
553 restrict the support of the learned policy by construction [19, 8, 7]. Our work follows this support-  
554 constrained perspective, but revisits the optimization problem that remains after support has been  
555 controlled. Namely, we identify that entropy regularization is essential to prevent policy collapse, and

556 leverage adjoint matching [16] to perform entropy-regularized policy extraction in latent space with  
 557 an *expressive* flow policy that achieves a better performance-entropy tradeoff than Gaussian policies.

558 **Entropy-regularized policy extraction.** The role of entropy regularization has been well identified  
 559 in *online* RL as a way to improve exploration [13, 12], policy robustness [13, 14] and smoothing  
 560 of the optimization landscape [15]. In *offline* RL, while some works note the benefits of entropy  
 561 regularization for exploration [35, 36], the importance of entropy regularization has been less  
 562 emphasized, perhaps due to the implicit entropy regularization effect from various behavior policy  
 563 regularization techniques. However, in our work where the support constraint is handled by the policy  
 564 parametrization and thus a separate behavior policy is not needed, the standalone role of entropy  
 565 regularization separate from support control becomes more crucial.

## 566 B Proof of Lemma 3.2

567 For a Lebesgue-measurable set  $X \subseteq \mathbb{R}^d$ , let  $\lambda(X)$  denote its Lebesgue measure. We first prove the  
 568 following technical lemma to express the entropy as a KL divergence and a constant.

569 **Lemma B.1.** *Fix a state  $s \in \mathcal{S}$ . For brevity, we omit the state dependence in the policy. Let*  
 570  *$L := \text{supp}(q_L(\cdot|s))$  denote the support of  $q_L(\cdot|s)$ . Suppose  $L$  has finite positive Lebesgue measure,*  
 571 *i.e.,  $0 < \lambda(L) < \infty$ , and  $\pi_L$  satisfies the latent-space support constraints (6), i.e.,*

$$\text{supp}(\pi_L) \subseteq \text{supp}(q_L(\cdot|s)). \quad (21)$$

572 Let  $q_L^{\text{unif}}$  be the uniform distribution over  $\text{supp}(q_L(\cdot|s))$ , i.e.,

$$q_L^{\text{unif}}(z) = \begin{cases} 1/\lambda(L) & z \in L \\ 0 & \text{otherwise} \end{cases}. \quad (22)$$

573 Then, the entropy of  $\pi_L$  can be expressed as

$$\mathcal{H}(\pi_L) = -D_{\text{KL}}(\pi_L \| q_L^{\text{unif}}) + \log \lambda(L), \quad (23)$$

574 where the second term in (23) is a constant independent of  $\pi_L$ .

575 *Proof.* Expanding the KL divergence,

$$D_{\text{KL}}(\pi_L \| q_L^{\text{unif}}) = \int_{\mathbb{R}^d} \pi_L(z) \log \frac{\pi_L(z)}{q_L^{\text{unif}}(z)} dz, \quad (24)$$

$$= \int_L \pi_L(z) \log \pi_L(z) dz - \int_L q_L^{\text{unif}}(z) \log q_L^{\text{unif}}(z) dz, \quad (25)$$

$$= -\mathcal{H}(\pi_L) - \int_L q_L^{\text{unif}}(z) \log\left(\frac{1}{\lambda(L)}\right) dz, \quad (26)$$

$$= -\mathcal{H}(\pi_L) + \log \lambda(L). \quad (27)$$

576 Rearranging the above equation yields the desired result.  $\square$

577 We are now ready to prove Lemma 3.2.

578 *Proof.* Fix a state  $s \in \mathcal{S}$ . For brevity, we omit the state dependence in the policy. Restating the  
 579 entropy-regularized optimization problem (11) in terms of the policy  $\pi_L$  directly, we wish to solve the  
 580 following optimization problem

$$\min_{\pi_L} -\mathbb{E}_{s \sim \mathcal{D}} [\mathbb{E}_{z \sim q_B} [Q^{\pi_L}(z)] + \alpha \mathcal{H}(\pi_L)]. \quad (28)$$

581 By Theorem B.1, the entropy of  $\pi_L$  can be expressed as

$$\mathcal{H}(\pi_L) = -D_{\text{KL}}(\pi_L \| q_L^{\text{unif}}) + \log \lambda(L). \quad (29)$$

582 Substituting this into (28) and dropping the constant term, we now have a KL divergence regularized  
 583 optimization problem

$$\min_{\pi_L} -\mathbb{E}_{s \sim \mathcal{D}} [\mathbb{E}_{z \sim q_B} [Q^{\pi_L}(z)] - \alpha D_{\text{KL}}(\pi_L \| q_L^{\text{unif}})]. \quad (30)$$

584 Now, define the partition function  $Z$  as

$$Z := \int_{\mathbb{R}^d} q_{\mathbb{L}}^{\text{unif}}(z) \exp\left(\frac{1}{\alpha} Q^{\pi_{\mathbb{L}}}(z)\right) dz. \quad (31)$$

585 Since  $Q^{\pi_{\mathbb{L}}}$  is measurable and bounded from above, there exists a  $m \in \mathbb{R}$  such that  $Q^{\pi_{\mathbb{L}}}(z) \leq m$  for  
586 all  $z \in \mathbb{R}^d$ . Hence,

$$0 < \exp\left(\frac{1}{\alpha} Q^{\pi_{\mathbb{L}}}(z)\right) \leq \exp\left(\frac{1}{\alpha} m\right) < \infty, \quad (32)$$

587 and thus

$$0 < Z \leq \exp\left(\frac{1}{\alpha} m\right) < \infty, \quad (33)$$

588 and  $Z$  is finite, thus  $\pi_{\mathbb{L}}^*$ , defined as

$$\pi_{\mathbb{L}}^*(z) := \frac{1}{Z} q_{\mathbb{L}}^{\text{unif}}(z) \exp\left(\frac{1}{\alpha} Q^{\pi_{\mathbb{L}}}(z)\right). \quad (34)$$

589 is well-defined. Since  $D_{\text{KL}}(\pi_{\mathbb{L}} \| q_{\mathbb{L}}^{\text{unif}}) = +\infty$  for any  $\pi_{\mathbb{L}} \not\ll q_{\mathbb{L}}^{\text{unif}}$ , it is sufficient to consider policies  
590  $\pi_{\mathbb{L}}$  that are absolutely continuous with respect to  $q_{\mathbb{L}}^{\text{unif}}$ . Note that

$$\frac{\pi_{\mathbb{L}}^*(z)}{q_{\mathbb{L}}^{\text{unif}}(z)} = \frac{1}{Z} \exp\left(\frac{1}{\alpha} Q^{\pi_{\mathbb{L}}}(z)\right), \quad \log\left(\frac{\pi_{\mathbb{L}}^*(z)}{q_{\mathbb{L}}^{\text{unif}}(z)}\right) = -\log Z + \frac{1}{\alpha} Q^{\pi_{\mathbb{L}}}(z). \quad (35)$$

591 For any  $\pi_{\mathbb{L}} \ll q_{\mathbb{L}}^{\text{unif}}$ , we have

$$D_{\text{KL}}(\pi_{\mathbb{L}} \| \pi_{\mathbb{L}}^*) = \int_{\mathbb{R}^d} \pi_{\mathbb{L}}(z) \log \frac{\pi_{\mathbb{L}}(z)}{\pi_{\mathbb{L}}^*(z)} dz, \quad (36)$$

$$= \int_{\mathbb{R}^d} \pi_{\mathbb{L}}(z) \log \frac{\pi_{\mathbb{L}}(z)}{q_{\mathbb{L}}^{\text{unif}}(z)} dz - \int_{\mathbb{R}^d} \pi_{\mathbb{L}}(z) \log \frac{\pi_{\mathbb{L}}^*(z)}{q_{\mathbb{L}}^{\text{unif}}(z)} dz, \quad (37)$$

$$= D_{\text{KL}}(\pi_{\mathbb{L}} \| q_{\mathbb{L}}^{\text{unif}}) + \log Z - \frac{1}{\alpha} \int_{\mathbb{R}^d} \pi_{\mathbb{L}}(z) Q^{\pi_{\mathbb{L}}}(z) dz, \quad (38)$$

$$= D_{\text{KL}}(\pi_{\mathbb{L}} \| q_{\mathbb{L}}^{\text{unif}}) + \log Z - \frac{1}{\alpha} \mathbb{E}_{z \sim \pi_{\mathbb{L}}}[Q^{\pi_{\mathbb{L}}}(z)]. \quad (39)$$

592 Thus, the objective can be written as the KL divergence to  $\pi_{\mathbb{L}}^*$  plus a constant:

$$-\mathbb{E}_{z \sim \pi_{\mathbb{L}}}[Q^{\pi_{\mathbb{L}}}(z) - \alpha D_{\text{KL}}(\pi_{\mathbb{L}} \| q_{\mathbb{L}}^{\text{unif}})] = \alpha D_{\text{KL}}(\pi_{\mathbb{L}} \| \pi_{\mathbb{L}}^*) - \alpha \log Z. \quad (40)$$

593 Since  $D_{\text{KL}}(\pi_{\mathbb{L}} \| \pi_{\mathbb{L}}^*) \geq 0$  for all  $\pi_{\mathbb{L}}$  with equality if and only if  $\pi_{\mathbb{L}} = \pi_{\mathbb{L}}^*$ , we conclude that  $\pi_{\mathbb{L}}^*$  is the  
594 unique optimal solution to the optimization problem (28).  $\square$

## 595 C More Theoretical Results

### 596 C.1 The Static Flow Issue and Implicit Entropy Regularization in ReFORM

597 To enforce the latent space support constraint (6), ReFORM [7] introduces a reflection term during  
598 the integration of the latent flow policy  $\mu_{\theta_{\mathbb{B}-\mathbb{L}}}$ . The corresponding modified Euler step is defined as:

$$w_{k+1} = \mathbb{1}\{\hat{w}_{k+1} \in \mathcal{B}_l^d\} \hat{w}_{k+1} + (1 - \mathbb{1}\{\hat{w}_{k+1} \in \mathcal{B}_l^d\}) (\hat{w}_{k+1} - \langle v_{\theta_{\mathbb{B}-\mathbb{L}}}(w_k, k\Delta t; s) \Delta t, n_{k+1} \rangle n_{k+1}), \quad (41)$$

599 where  $\hat{w}_{k+1} = w_k + v_{\theta_{\mathbb{B}-\mathbb{L}}}(k\Delta t, w_k; s) \Delta t$  denotes the standard Euler step,  $N$  is the total number  
600 of integration steps,  $k \in \{0, \dots, N-1\}$ ,  $\Delta t = \frac{1}{N}$ ,  $n_{k+1} = \frac{\hat{w}_{k+1}}{\|\hat{w}_{k+1}\|}$ , and  $\langle \cdot, \cdot \rangle$  represents the inner  
601 product. Integrating this modified step yields the final latent variable  $z \leftarrow w_N$ .

602 While [7] proves that this modified integration step guarantees adherence to the latent support  
603 constraint (6), it introduces a degenerate behavior we term ‘‘static flow’’. This occurs when the  
604 proposed step violates the boundary ( $\hat{w}_{k+1} \notin \mathcal{B}_l^d$ ) and the flow velocity is strictly codirectional with  
605 the outward normal vector ( $v_{\theta_{\mathbb{B}-\mathbb{L}}} \uparrow \uparrow n_{k+1}$ ). Under these conditions, the reflection term perfectly  
606 cancels the velocity step:

$$\begin{aligned} w_{k+1} &= \hat{w}_{k+1} - \langle v_{\theta_{\mathbb{B}-\mathbb{L}}}(w_k, k\Delta t; s) \Delta t, n_{k+1} \rangle n_{k+1} \\ &= (w_k + v_{\theta_{\mathbb{B}-\mathbb{L}}}(w_k, k\Delta t; s) \Delta t) - v_{\theta_{\mathbb{B}-\mathbb{L}}}(w_k, k\Delta t; s) \Delta t = w_k. \end{aligned} \quad (42)$$

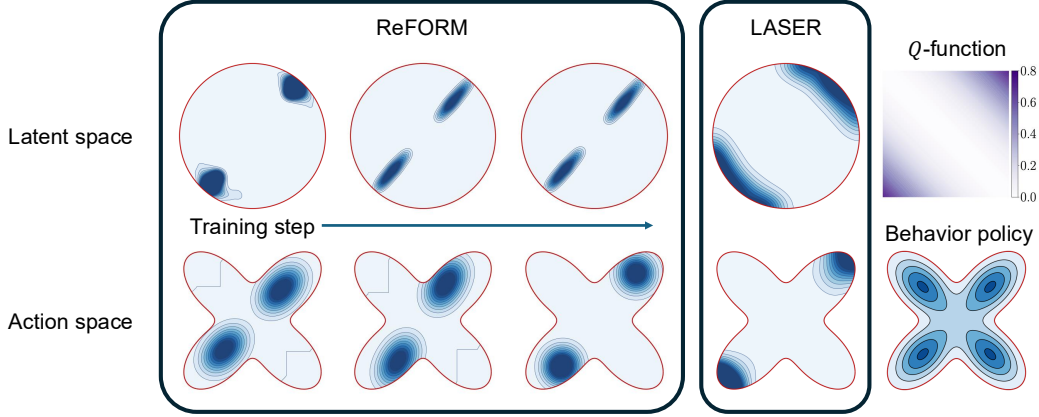


Figure 7: **ReFORM fails to maximize the  $Q$ -value because of the “static flow” issue.** For a given behavior policy and a ground-truth  $Q$ -function, the generated latent space distribution  $(\mu_{\theta_{B \rightarrow L}}^s)_{\#} q_B$  concentrates on the radius that has the same direction as the flow, instead of moving towards the boundary of  $\mathcal{B}_l^d$ . This makes the flow static, resulting in suboptimal performance compared with **LASER**.

607 Consequently, the trajectory becomes stationary. Rather than following the flow field  $v_{\theta_{B \rightarrow L}}$  to the  
 608 boundary of  $\mathcal{B}_l^d$ , the pushforward distribution of  $\mu_{\theta_{B \rightarrow L}}$  unnaturally accumulates along the radial  
 609 trajectory. To illustrate this phenomenon, we train ReFORM on the 2D toy environment introduced in  
 610 Figure 3. The learned latent distribution and its resulting action distribution are visualized in Figure 7.  
 611 We empirically observe that the latent mass inappropriately concentrates along two radii, resulting in  
 612 a suboptimal action distribution.

613 Interestingly, we find that this static flow artifact unintentionally serves as an implicit entropy  
 614 regularizer, which partially explains ReFORM’s empirical success. As demonstrated in Figure 7, the  
 615 premature halting of trajectories prevents the latent distribution from collapsing to a single point,  
 616 artificially maintaining the spread of the distribution.

## 617 C.2 The Ground-truth Base Velocity Field from the Base Space to the Latent Space

618 As mentioned in Section 3.2, there exists a ground-truth velocity field  $v_{B \rightarrow L}$  that maps samples from  
 619  $w \sim q_B = \mathcal{N}(0, \sigma_0 I)$  to  $z \sim \mathcal{U}(\mathcal{B}_l^d)$ . We show here how we can find such a ground-truth velocity  
 620 field and why we choose not to use it in our implementation.

621 Since both the base distribution  $q_B = \mathcal{N}(0, \sigma_0^2 I)$  and the target distribution  $\mathcal{U}(\mathcal{B}_l^d)$  are spherically  
 622 symmetric, the optimal transport (OT) mapping between them operates purely along the radial  
 623 direction. Let  $r_w = \|w\|$  denote the radius of a sample  $w$  from the base distribution, and  $r_z = \|z\|$   
 624 denote the radius in the latent space. The cumulative distribution function (CDF) for the uniform  
 625 distribution on a  $d$ -dimensional ball of radius  $l$  is  $F_{\mathcal{U}}(r) = (r/l)^d$  for  $r \in [0, l]$ . Meanwhile, the  
 626 radius of the Gaussian base distribution follows a scaled  $\chi$ -distribution with  $d$  degrees of freedom,  
 627 with its CDF denoted as  $F_{\chi}(r)$ .

628 By matching the quantiles of these radial distributions, we can construct the exact analytical mapping  
 629  $\mu_{B \rightarrow L} : \mathbb{R}^d \rightarrow \mathcal{B}_l^d$  as follows:

$$\mu_{B \rightarrow L}(w) = \frac{w}{\|w\|} F_{\mathcal{U}}^{-1}(F_{\chi}(\|w\|)) = w \frac{l \cdot F_{\chi}(\|w\|)^{1/d}}{\|w\|}. \quad (43)$$

630 However, although we have derived an exact analytical mapping  $\mu_{B \rightarrow L}$ , it is nontrivial to derive the  
 631 corresponding velocity field  $v_{B \rightarrow L}(w_t, t; s)$ , i.e., the instantaneous velocity of a particle currently  
 632 at location  $w_t$  at time  $t$ . Under standard optimal transport displacement interpolation, a particle  
 633 travels in a straight line from its origin  $w$  to its destination  $\mu_{B \rightarrow L}(w)$  at a constant speed, yielding the  
 634 trajectory:

$$w_t = (1 - t)w + t\mu_{B \rightarrow L}(w). \quad (44)$$

635 While the velocity of this specific particle is trivially  $v(w) = \mu_{\mathcal{B} \rightarrow \mathcal{L}}(w) - w$ , the velocity field must be  
 636 expressed independently of the origin  $w$ . This requires analytically inverting the trajectory equation  
 637 to express  $w$  purely as a function of the current state  $w_t$  and time  $t$ .

638 Because our mapping  $\mu_{\mathcal{B} \rightarrow \mathcal{L}}$  is purely radial, the vectors  $w$  and  $w_t$  share the same direction, meaning  
 639  $\frac{w}{\|w\|} = \frac{z_t}{\|z_t\|}$ . Let  $r_w = \|w\|$  and  $r_t = \|w_t\|$ . By taking the norm of the trajectory equation, the  
 640 relationship between the magnitudes simplifies to:

$$r_t = (1 - t)r_w + t \cdot l \cdot F_\chi(r_w)^{1/d}. \quad (45)$$

641 To compute the velocity at  $w_t$ , we must find the inverse function  $g_t(r_t)$  that recovers the initial radius  
 642  $r_w$  given  $r_t$  and  $t$ . If such an inverse existed, the exact Eulerian velocity field would be:

$$v_{\mathcal{B} \rightarrow \mathcal{L}}(w_t, t) = \frac{w_t - w}{t} = \left( \frac{\|w_t\| - g_t(\|w_t\|)}{t} \right) \frac{w_t}{\|w_t\|}. \quad (46)$$

643 However, the CDF of the scaled  $\chi$ -distribution  $F_\chi$  involves the regularized lower incomplete gamma  
 644 function. Consequently, the trajectory equation governing  $r_t$  and  $r_w$  possesses no closed-form  
 645 algebraic inverse  $g_t$ . To evaluate  $v_{\mathcal{B} \rightarrow \mathcal{L}}(t, z)$  accurately at runtime, one would have to execute a  
 646 numerical root-finding algorithm to approximate  $g_t(\|z\|)$  at every single evaluation step of the ODE  
 647 solver. This requirement is computationally prohibitive and highly susceptible to numerical instability,  
 648 particularly in high dimensions where the  $\chi$ -distribution concentrates sharply, or near the boundaries  
 649 where gradients vanish. By parameterizing the velocity field with a neural network  $v_{\theta_{\mathcal{B} \rightarrow \mathcal{L}}}$  and learning  
 650 it via flow matching, we completely sidestep this intractable analytical inversion, yielding an efficient  
 651 and stable model for inference.

## 652 D Practical Implementation of LASER

653 We provide a high level summary of LASER in Algorithm 1.

654 **Flow-matching Decoder Learning.** We parameterize the velocity field  $v_{\theta_{\mathcal{L} \rightarrow \mathcal{A}}}$  of the decoder  $\mu_{\theta_{\mathcal{L} \rightarrow \mathcal{A}}}$   
 655 with multilayer perceptrons (MLP), and learn it via the flow matching loss (8). The prior latent  
 656 space distribution is chosen to be a bounded uniform distribution over a  $d$ -ball with radius  $l$ , i.e.,  
 657  $q_{\mathcal{L}} = \mathcal{U}(\mathcal{B}_l^d)$ . After training, we sample  $z \sim q_{\mathcal{L}}$  and use Euler integration to get the action samples:

$$z_0 \leftarrow z, \quad z_{t+\Delta t} = z_t + v_{\theta_{\mathcal{L} \rightarrow \mathcal{A}}}(z_t, k\Delta t; s)\Delta t, \quad (47)$$

658 where  $N$  is the total number of integration steps,  $\Delta t = \frac{1}{N}$ ,  $k \in \{0, \dots, N - 1\}$ , and the final  
 659 generated action is  $a \leftarrow z_1$ .

660 **Flow-matching for  $q_{\mathcal{L}}^{\text{unif}}$ .** Similarly, the velocity field  $v_{\mathcal{B} \rightarrow \mathcal{L}}^{\text{base}}$  is parameterized with MLP and learned  
 661 via the flow matching loss (8), with prior distribution  $q_{\mathcal{B}} = \mathcal{N}(0, \sigma_0^2 I)$  and the desired pushforward  
 662  $q_{\mathcal{L}}^{\text{unif}} = \mathcal{U}(\mathcal{B}_l^d)$ . Note that the  $v_{\mathcal{B} \rightarrow \mathcal{L}}^{\text{base}}$  is state-independent.

663 **Critic Learning.** The critic is learned via the standard TD learning objective (2) with  $K$  ensembles.  
 664 In some environments, we found that using a pessimistic target-value backup [66] can improve  
 665 performance. In that case, the loss function becomes

$$L(\phi_j) = \mathbb{E}_{(s, a, s') \sim \mathcal{D}, a' \sim \pi_{\mathcal{A}}(\cdot | s')} \left[ \left( Q_{\phi_j}(s, a) - r - \gamma \left( Q_{\hat{\phi}}^{\text{mean}}(s', a') - \rho Q_{\hat{\phi}}^{\text{std}}(s', a') \right) \right)^2 \right], \quad (48)$$

666 where  $\phi_j$  is the parameters for the  $j$ -th ensemble,  $Q_{\hat{\phi}}^{\text{mean}}(s', a') := \frac{1}{K} \sum_k Q_{\hat{\phi}_k}(s', a')$ ,  $Q_{\hat{\phi}}^{\text{std}}(s', a') =$   
 667  $\sqrt{\sum_k \left( Q_{\hat{\phi}_k}(s', a') - Q_{\hat{\phi}}^{\text{mean}}(s', a') \right)^2}$ , and  $\rho$  is the pessimistic coefficient.

668 **Entropy-regularized Policy Optimization via Adjoint Matching.** Practically, we solve all the  
 669 ODEs in the adjoint matching process by Euler integration with fixed step size  $\Delta t = 1/N$ , where  $N$   
 670 is the number of discretization steps. First, we integrate the forward SDE (16) with

$$w_{t+\Delta t} \leftarrow w_t + \left( 2v_{\theta_{\mathcal{B} \rightarrow \mathcal{L}}}(w_t, t; s) - \frac{w_t}{t} \right) \Delta t + \sigma_0 \sqrt{\frac{2\Delta t(1-t)}{t}} n_t, \quad w_0 \sim \mathcal{N}(0, \sigma_0^2 I), \quad (49)$$

---

**Algorithm 1** LASER
 

---

```

1: while not converged do
2:   ▷ Train Critic
3:   Sample batch  $\{(s, a, r, s')\} \sim \mathcal{D}$ 
4:    $w \sim \mathcal{N}(0, \sigma_0^2 I_d)$ 
5:    $a' = \mu_{\theta_{L \rightarrow A}}^{s'}(\mu_{\theta_{B \rightarrow L}}^{s'}(w))$ 
6:   Update  $\phi$  to minimize  $L(\phi) = \mathbb{E}[(r(s, a) + \gamma Q_{\hat{\phi}}^{\pi_A}(s', a') - Q_{\phi}^{\pi_A}(s, a))^2]$  (2)

7:   ▷ Train vector field  $v_{\theta_{L \rightarrow A}}$  in decoder  $\mu_{\theta_{L \rightarrow A}}$ 
8:    $x^0 \sim q_L$ 
9:    $x^1 \leftarrow a$ 
10:   $t \sim \mathcal{U}([0, 1])$ 
11:   $x^t \leftarrow (1 - t)x^0 + tx^1$ 
12:  Update  $\theta_{L \rightarrow A}$  to minimize  $L_{L \rightarrow A}(\theta_{L \rightarrow A}) = \mathbb{E}[\|v_{\theta_{L \rightarrow A}}(x_t, t; s) - (a - z)\|^2]$  (8)

13:  ▷ Train vector field  $v_{B \rightarrow L}^{\text{base}}$  to match  $q_L^{\text{unif}}$ 
14:   $x^0 \sim q_B$ 
15:   $x^1 \sim q_L^{\text{unif}}$ 
16:   $t \sim \mathcal{U}([0, 1])$ 
17:   $x^t \leftarrow (1 - t)x^0 + tx^1$ 
18:  Update  $v_{B \rightarrow L}^{\text{base}}$  to minimize  $\mathbb{E}[\|v_{B \rightarrow L}^{\text{base}}(x_t, t; s) - (a - z)\|^2]$ 

19:  ▷ Train Latent-Space Policy using Adjoint Matching
20:   $w_0 \sim \mathcal{N}(0, \sigma_0^2 I)$ 
21:   $w_{0+\Delta t}, \dots, w_1 \leftarrow$  Forward simulate SDE (16) with (49)
22:   $\tilde{g}_1 \leftarrow -\frac{1}{\alpha} \nabla_{w_1} Q_{\phi}^{\pi_L}(s, w_1)$ 
23:   $\tilde{g}_{1-\Delta t}, \dots, \tilde{g}_0 \leftarrow$  solve ODE (19) backwards
24:  Compute regression targets  $\hat{v}_t^*$  for  $t = 0, \dots, 1$  using (18)
25:  Update  $v_{\theta_{B \rightarrow L}}$  to minimize  $L_{\text{AM}}(\theta_{B \rightarrow L}) = \mathbb{E} \left[ \frac{1}{2N} \sum_t \left\| \frac{2}{\sigma_t} \left( v_{\theta_{B \rightarrow L}}(w_t, t; s) - \hat{v}_t^* \right) \right\|^2 \right]$  (20)
26: return policy  $\pi_A$ 

```

---

671 where  $n_t \sim \mathcal{N}(0, I)$ . Then, we calculate the “lean” adjoint states by integrating the reverse ODE  
 672 from  $\tilde{g}_1 = -\frac{1}{\alpha} \nabla_{w_1} Q_{\phi}^{\pi_L}(s, w_1)$  following:

$$\tilde{g}_{t-\Delta t} \leftarrow \tilde{g}_t + \text{VJP} \left( \nabla_{w_t} \left( 2v_{B \rightarrow L}^{\text{base}}(w_t, t; s) - \frac{w_t}{t} \right), \tilde{g}_t \right) \Delta t, \quad (50)$$

673 where  $\text{VJP}(\cdot, \cdot)$  is the vector-Jacobian product. Then, we can calculate the optimal control  $\hat{v}^*$   
 674 following (18) and then calculate and optimize loss (20). Following [16], we also apply reward  
 675 clipping to stabilize training, i.e., after calculating  $\tilde{g}_1$ , we clip the top 10% largest values to the 10%  
 676 threshold.

## 677 E Details of the Reach-avoid Example in Figure 1

678 In this environment, the agents start from two separate starting ranges at the bottom and tend to reach  
 679 the top goal. The action space is the 2-dimensional velocity. When calculating the agents’ next state,  
 680 if the agents’ next state is inside the obstacles, we set  $s' = s$ . The “No Entropy Regularization”  
 681 baseline is constructed by replacing the adjoint matching component in LASER with plain gradient  
 682 ascent, and the generated latent samples  $z$  are squashed inside the ball  $\mathcal{B}_t^d$  to ensure support constraint.  
 683 The “No Support Constraint” baseline is constructed with a QAM [23] model with a large inverse  
 684 temperature  $\tau$  that drives the action distribution away from the behavior policy distribution.

## 685 F Experiments

### 686 F.1 Environments

687 We conduct experiments on the widely used OGBench offline RL benchmark [18]. Although  
688 OGBench is designed for offline goal-conditioned RL, it provides the `-singletask` flag for standard  
689 offline RL. We use 4 environments, including one locomotion environment (`antmaze`) and three  
690 manipulation environments (`cube-single`, `cube-double`, and `scene`), each with 5 tasks. For each task,  
691 we consider 2 datasets, yielding a total of 40 tasks.

692 We consider the following tasks whose given dataset is `CLEAN`, i.e., the demonstrations are randomly  
693 generated by an expert policy:

- 694 • `antmaze-large-navigate-singletask-task{1, 2, 3, 4, 5}-v0`
- 695 • `cube-single-play-singletask-task{1, 2, 3, 4, 5}-v0`
- 696 • `cube-double-play-singletask-task{1, 2, 3, 4, 5}-v0`
- 697 • `scene-play-singletask-task{1, 2, 3, 4, 5}-v0`

698 In addition, we consider the following tasks with the `NOISY` dataset, i.e., the demonstrations are  
699 generated by a highly noisy policy:

- 700 • `antmaze-large-explore-singletask-task{1, 2, 3, 4, 5}-v0`
- 701 • `cube-single-noisy-singletask-task{1, 2, 3, 4, 5}-v0`
- 702 • `cube-double-noisy-singletask-task{1, 2, 3, 4, 5}-v0`
- 703 • `scene-noisy-singletask-task{1, 2, 3, 4, 5}-v0`

704 More details of the environments and datasets can be found in [18].

### 705 F.2 Baselines

706 In this section, we briefly introduce the details of each baseline.

707 **FQL.** FQL [21] is an algorithm that uses the statistical distance regularization objective (4), instanti-  
708 ating  $D(\pi_\theta \parallel \pi_\beta)$  as the Wasserstein distance. It begins by training the decoder  $\mu_{\theta_{L \rightarrow A}}^s$  via the flow  
709 matching loss (8), assuming a standard normal latent distribution  $q_L = \mathcal{N}(0, I)$ . Subsequently, it  
710 fine-tunes this decoder by learning a *one-step* distilled policy  $\mu_{\theta'_{L \rightarrow A}}(z, s) : \mathcal{L} \times \mathcal{S} \rightarrow \mathcal{A}$ . This policy  
711 is optimized to maximize expected return while remaining anchored to the behavior cloning (BC)  
712 decoder  $\mu_{\theta_{L \rightarrow A}}^s$ . Specifically, the one-step policy minimizes the following loss:

$$L(\theta'_{L \rightarrow A}) = \mathbb{E}_{s \sim D, z \sim q_L} \left[ -Q_\phi^{\mu_{\theta'_{L \rightarrow A}}}(s, \mu_{\theta'_{L \rightarrow A}}(z, s)) \right] + \alpha \mathbb{E}_{s \sim D, z \sim q_L} \left[ \|\mu_{\theta'_{L \rightarrow A}}(z, s) - \mu_{\theta_{L \rightarrow A}}^s(z)\| \right], \quad (51)$$

713 where the first term optimizes task performance, and the second term enforces the statistical distance  
714 regularization. Although the one-step policy is less expressive than the original flow policy [22],  
715 training this policy to maximize the  $Q$ -value avoids the BPTT process. However, the Wasserstein  
716 distance regularization is not enough to avoid OOD actions [7].

717 **QAM.** QAM [23] is an algorithm that uses the statistical distance regularization objective (4), instanti-  
718 ating  $D(\pi_\theta \parallel \pi_\beta)$  as the KL divergence. It begins by training the decoder  $\mu_{\theta_{L \rightarrow A}}^s$  via the flow matching  
719 loss (8), assuming a standard normal latent distribution  $q_L = \mathcal{N}(0, I)$ . Subsequently, it fine-tunes the  
720 decoder via adjoint matching to align with the target tilted distribution  $\pi_\beta(\cdot|s) \exp(\tau Q^{\mu_{\theta_{L \rightarrow A}}^s}(s, \cdot))$ .  
721 This procedure is mathematically equivalent to optimizing the following KL-divergence-regularized  
722 objective:

$$L(\theta) = -\mathbb{E}_{s \sim D} \left[ \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [Q^{\mu_{\theta_{L \rightarrow A}}^s}(s, a)] + \frac{1}{\tau} D_{\text{KL}}(\pi_\theta(\cdot|s) \parallel \pi_\beta(\cdot|s)) \right], \quad (52)$$

723 where  $\pi_\theta(\cdot|s) = (\mu_{\theta_{L \rightarrow A}}^s)_\# q_L$ , and  $\mu_{\theta'_{L \rightarrow A}}^s$  is the fine-tuned flow decoder. QAM avoids the BPTT process  
724 by leveraging the adjoint matching technique, similar to `LASER`, but the KL divergence regularization  
725 can limit the policy improvement [7].

726 **QAM-E.** QAM-E [23] is an improved version on QAM, which after training the fine-tuned decoder  
 727  $\mu_{\theta_{L \rightarrow A}}^s$ , they learn another edit policy  $\pi_{\omega}(a|s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{A})$  to modify the action output from  
 728  $\pi_{\theta}(\cdot|s) = (\mu_{\theta_{L \rightarrow A}}^s)_{\#} q_L$  by at most  $\sigma_a$ , to additionally maximize the performance. In practice, the  
 729 edit policy  $\pi_{\omega}$  is parameterized as a squashed Gaussian policy that only outputs values in  $[-\sigma_a, \sigma_a]$ .  
 730 Empirically, QAM-E outperforms QAM [23], but it is sensitive to  $\sigma_a$ , and the edit policy distribution  
 731 may generate OOD actions.

732 **IFQL.** Presented in [21], IFQL is the flow version of the implicit diffusion Q-learning (IDQL)  
 733 algorithm [32]. Similarly to the aforementioned algorithms, it first learns the decoder  $\mu_{\theta_{L \rightarrow A}}^s$  via the  
 734 flow matching loss (8), assuming  $q_L = \mathcal{N}(0, I)$ . However, instead of learning the policy  $Q$ -function  
 735 with an actor-critic structure with loss (2), they learn the approximately optimal critic via implicit  
 736 Q-learning (IQL) [25] with the following loss:

$$L(\phi) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} [(Q_{\phi}(s, a) - r - V_{\xi}(s'))^2], \quad (53)$$

$$L(\xi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^{\kappa}(Q_{\hat{\phi}}(s, a) - V_{\xi}(s))], \quad (54)$$

737 where  $L_2^{\kappa}(u) = |\kappa - \mathbb{1}(u < 0)|u^2$  is the expectile regression loss. To extract a policy from the  
 738 learned  $Q$ -function and the decoder, IFQL uses rejection sampling:

$$a^* \leftarrow \arg \max_{a_1, \dots, a_N} Q(s, a_i), \quad a_1, \dots, a_N \sim (\mu_{\theta_{L \rightarrow A}}^s)_{\#} q_L. \quad (55)$$

739 IFQL satisfies the support constraint by construction, but the rejection sampling is inefficient in  
 740 high-dimensional spaces with a complex  $Q$ -function.

741 **DSRL.** DSRL [8] is a recently proposed latent space RL algorithm, which has several key differences  
 742 compared with LASER. First, the latent distribution of the decoder is a standard Gaussian in DSRL, i.e.,  
 743  $q_L = \mathcal{N}(0, I)$ . Because the support of  $q_L$  is unbounded, steering  $q_L$  might generate a distribution that  
 744 has high density far from the origin, where the samples are not well represented during the training of  
 745 the decoder. To avoid this problem, DSRL parameterizes the latent policy  $\pi_L$  with a squashed Gaussian  
 746 distribution that only generates samples within the range of  $[-\sigma_z, \sigma_z]$ , where  $\sigma_z$  is a hyperparameter  
 747 that requires fine-tuning for each task. DSRL follows a SAC-style loss [12] and explicitly regularizes  
 748 the entropy, which is straightforward for a Gaussian latent policy. However, because of the Gaussian  
 749 parameterization, DSRL has a suboptimal entropy-performance trade-off as shown in Figure 3.

750 **ReFORM.** ReFORM [7] is another recently proposed latent space RL algorithm, which also uses a  
 751 flow latent policy. The BC decoder  $\mu_{\theta_{L \rightarrow A}}^s$  learning process of ReFORM is the same as LASER. However,  
 752 unlike LASER, ReFORM does not train the latent policy via adjoint matching, but via plain gradient  
 753 descent. In addition, ReFORM does not have the explicit entropy regularization term, resulting in the  
 754 following loss:

$$L_{B \rightarrow L}(\theta_{B \rightarrow L}) = -\mathbb{E}_{s \sim \mathcal{D}, w \sim q_B} [Q^{\pi_L}(s, \mu_{\theta_{B \rightarrow L}}(w; s))]. \quad (56)$$

755 Directly using gradient descent on this loss cannot guarantee the latent space support constraint (6),  
 756 as the flow  $\mu_{\theta_{B \rightarrow L}}$  can drive the samples outside the closed ball  $\mathcal{B}_r^d$ . To address this issue, ReFORM  
 757 modifies the Euler integration step to (41) when integrating the flow  $\mu_{\theta_{B \rightarrow L}}$ , and proves that this  
 758 modified step guarantees latent space support constraint. However, as we introduced in Appendix C.1,  
 759 this modified step causes a “static flow” issue, which acts as an implicit entropy regularization but  
 760 makes the final action distribution suboptimal.

### 761 F.3 Hyperparameters

762 We organize the hyperparameters into three distinct categories. **Fixed common hyperparameters**  
 763 (Table 1) are shared across all algorithms and remain constant across all tasks. We do not fine-tune  
 764 these, but instead adopt the standard values established in prior work [21, 23, 7]. **Environment-**  
 765 **specific common hyperparameters** (Table 2) are also shared across algorithms but vary by task to  
 766 optimize baseline performance. Following [23], we perform a minimal sweep over two candidate  
 767 values for these settings: 1 or 5 for the action chunking length, and 0 or 0.5 for the pessimistic  
 768 backup weight. **Environment-specific unique hyperparameters** (Table 3) are specific to individual  
 769 baselines (namely FQL, DSRL, QAM, and QAM-E) and require per-environment tuning. For these, we

Table 1: **Fixed common hyperparameters.** These hyperparameters are shared by [LASER](#) and all the baselines, and are fixed in all environments.

Hyperparameter	Value
Learning rate	0.0003
Optimizer	Adam [67]
Maximum gradient norm	1 ( <a href="#">QAM</a> , <a href="#">QAM-E</a> ), 10 (others)
Target network smoothing coefficient	0.005
Discount factor $\gamma$	0.995
MLP dimensions	[512, 512, 512, 512]
Nonlinearity	GELU [68]
Flow integration steps	10
Flow time sampling distribution	$\mathcal{U}[0, 1]$
Minibatch size	256
Number of ensembles for the $Q$ -function	8
$Q$ -function aggregation method	mean
Actor layer normalization [69]	False
Critic layer normalization [69]	True

Table 2: **Environment-specific common hyperparameters.** These hyperparameters are shared by [LASER](#) and all baselines, but are different per environment and dataset.

Environment	Dataset	Action chunk length	Pessimistic backup weight $\rho$
antmaze-large-navigate	CLEAN	1	0.5
antmaze-large-explore	NOISY	1	0.5
cube-single-play	CLEAN	5	0.5
cube-single-noisy	NOISY	1	0
cube-double-play	CLEAN	5	0.5
cube-double-noisy	NOISY	1	0
scene-play	CLEAN	5	0.5
scene-noisy	NOISY	5	0.5

770 rely on the optimal values reported in their respective papers or tune them empirically. One notable  
 771 exception is the edit scale  $\sigma_a$  for [QAM-E](#). While the original authors sometimes set  $\sigma_a = 0$  [23] that  
 772 renders it identical to [QAM](#), we enforce a minimum value of  $\sigma_a = 0.01$  to ensure the algorithms  
 773 remain distinct. Finally, while [LASER](#) introduces two unique hyperparameters, we keep them strictly  
 774 fixed across all tasks at  $1/\alpha = 10$  and  $p(\sigma_0) = 0.999$ .

#### 775 **F.4 Computation Resources**

776 Most of our experiments are run on a 13th Gen Intel(R) Core(TM) i7-13700KF CPU with 64GB  
 777 RAM and an NVIDIA GeForce RTX 4090 GPU. We report the training time of all the algorithms in  
 778 Table 4.

#### 779 **F.5 Additional Results**

780 **Ablation study of the standard deviation of the base distribution.** We study the sensitivity  
 781 of [LASER](#) w.r.t. the standard deviation  $\sigma_0$  of  $q_B$ . We vary  $\sigma_0$  such that the ball  $\mathcal{B}_l^d$  contains the  
 782  $p$ -confidence level of the Gaussian distribution  $\mathcal{N}(0, \sigma_0^2 I)$ . We report the results in Figure 8. The  
 783 results show that [LASER](#) is not sensitive to  $\sigma_0$ .

784 **Ablation study of the reward clipping technique.** A known drawback of adjoint matching is its  
 785 tendency to produce large gradients, which can destabilize the training process [16, 23]. To address  
 786 this, we employ reward clipping (mentioned in Appendix D), following [16]. Notably, [QAM](#) encounters  
 787 a similar challenge and mitigates it using gradient clipping instead. In Figure 8, we compare the  
 788 training curves across three configurations: an unclipped baseline, reward clipping, and gradient  
 789 clipping (`max_grad_norm = 1`). The results demonstrate that while the unclipped baseline suffers

Table 3: **Environment-specific unique hyperparameters.** These hyperparameters are unique for some algorithms and require fine-tuning per environment. [LASER](#), [ReFORM](#), and [IFQL](#) do not have such hyperparameters, thus do not require environment-specific fine-tuning.

Environment	Dataset	FQL ( $\alpha$ )	DSRL ( $\sigma_z$ )	QAM ( $\tau$ )	QAM-E ( $\tau, \sigma_a$ )
antmaze-large-navigate	CLEAN	3.0	1.25	10.0	(1.0, 0.1)
antmaze-large-explore	NOISY	0.3	1.25	10.0	(1.0, 0.1)
cube-single-play	CLEAN	300	0.5	1.0	(1.0, 0.01)
cube-single-noisy	NOISY	30	0.5	1.0	(3.0, 0.5)
cube-double-play	CLEAN	300	1.5	1.0	(1.0, 0.01)
cube-double-noisy	NOISY	30	1.5	1.0	(1.0, 0.01)
scene-play	CLEAN	300	0.75	1.0	(1.0, 0.01)
scene-noisy	NOISY	30	0.75	1.0	(1.0, 0.01)

Table 4: **Training time comparison.** We report the number of iterations per second (it/s) for [LASER](#) and all the baselines.

Algorithm	FQL	QAM	QAM-E	IFQL	DSRL	ReFORM	LASER
it/s	220	163	125	210	165	138	123

790 from severe instability, both reward clipping and gradient clipping successfully resolve the issue and  
 791 yield comparable overall performance.

792 **Detailed results.** The full results for each algorithm across all tasks are provided in Table 5. Success  
 793 rates are bolded for algorithms achieving at least 95% of the best performance per task. We observe  
 794 that [LASER](#) yields the best overall performance, approaching near-perfect success rates on multiple  
 795 tasks.

796 **Training curves.** All algorithms’ training curves of all tasks with both the CLEAN and the NOISY  
 797 are provided in Figure 9 - Figure 16.

## 798 F.6 Code

799 We provide the code of [LASER](#) in ‘laser.zip’ in the supplementary materials.

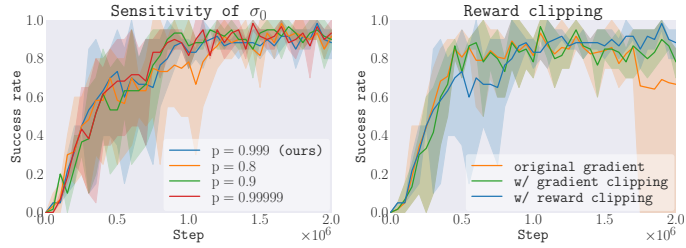


Figure 8: **Additional ablation studies.** LASER demonstrates robustness to the choice of  $\sigma_0$ . Furthermore, applying reward clipping effectively improves training stability, with gradient clipping yielding a comparable stabilizing effect.

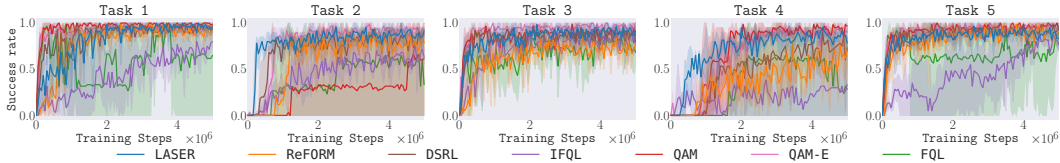


Figure 9: Training curves for the antmaze-large environment with the CLEAN dataset.

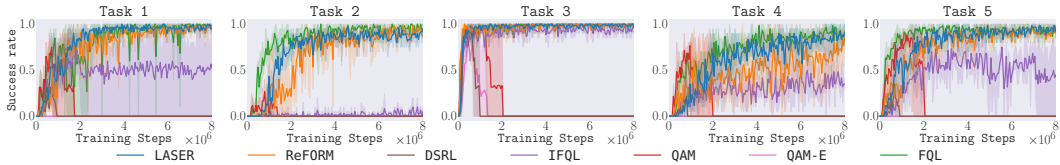


Figure 10: Training curves for the antmaze-large environment with the NOISY dataset.

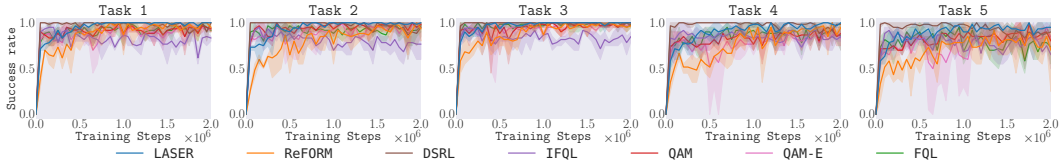


Figure 11: Training curves for the cube-single environment with the CLEAN dataset.

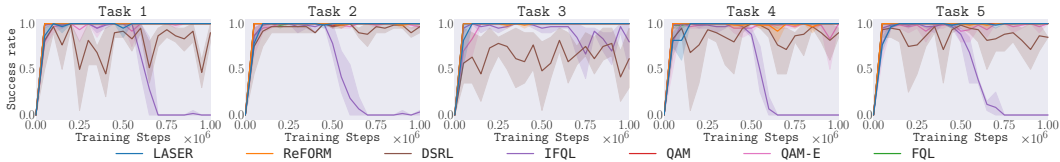


Figure 12: Training curves for the cube-single environment with the NOISY dataset.

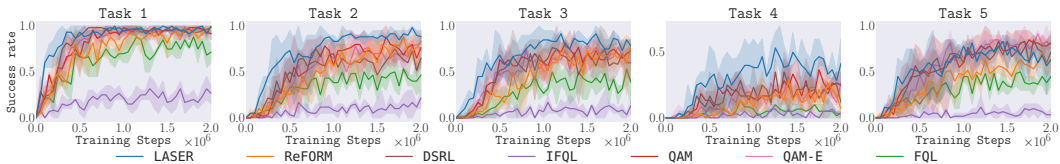


Figure 13: Training curves for the cube-double environment with the CLEAN dataset.

Table 5: **Full results on 40 OGBench tasks.** The results are over 3 seeds and 50 runs per seed. The results are bolded for algorithms achieves at or above 95% of the best performance following [18]. We omit the -singletask tags from the task names to save space.

Task	Dataset	FQL	QAM	QAM-E	IFQL	DSRL	ReFORM	LASER
antmaze-large-navigate-task1-v0	CLEAN	64±45	<b>96±3</b>	<b>93±2</b>	79±13	90±3	91±2	<b>95±1</b>
antmaze-large-navigate-task2-v0	CLEAN	31±43	60±42	<b>94±3</b>	65±5	81±5	78±2	<b>89±2</b>
antmaze-large-navigate-task3-v0	CLEAN	71±12	85±3	<b>97±1</b>	85±9	<b>94±3</b>	74±6	<b>92±3</b>
antmaze-large-navigate-task4-v0	CLEAN	33±41	<b>93±1</b>	<b>91±1</b>	31±43	75±14	69±4	85±1
antmaze-large-navigate-task5-v0	CLEAN	88±6	<b>95±1</b>	<b>93±4</b>	83±2	<b>91±4</b>	83±1	<b>91±1</b>
antmaze-large-explore-task1-v0	NOISY	<b>97±1</b>	0±0	0±0	50±36	0±0	91±4	<b>97±2</b>
antmaze-large-explore-task2-v0	NOISY	<b>97±1</b>	0±0	0±0	3±2	0±0	88±9	87±2
antmaze-large-explore-task3-v0	NOISY	<b>100±0</b>	0±0	0±0	<b>95±1</b>	0±0	<b>100±0</b>	<b>98±3</b>
antmaze-large-explore-task4-v0	NOISY	<b>89±3</b>	0±0	0±0	36±13	0±0	82±9	<b>89±2</b>
antmaze-large-explore-task5-v0	NOISY	<b>94±6</b>	0±0	0±0	47±34	0±0	<b>93±2</b>	<b>92±3</b>
cube-single-play-task1-v0	CLEAN	<b>95±4</b>	94±3	<b>95±4</b>	79±2	95±5	<b>99±2</b>	<b>100±0</b>
cube-single-play-task2-v0	CLEAN	<b>97±1</b>	<b>99±1</b>	<b>99±2</b>	70±5	90±7	<b>96±3</b>	<b>100±0</b>
cube-single-play-task3-v0	CLEAN	<b>99±1</b>	<b>100±0</b>	<b>99±1</b>	82±7	<b>97±4</b>	<b>98±3</b>	<b>100±0</b>
cube-single-play-task4-v0	CLEAN	90±4	91±2	85±4	82±3	<b>98±2</b>	87±2	<b>99±1</b>
cube-single-play-task5-v0	CLEAN	81±2	<b>89±4</b>	85±4	75±6	85±3	86±13	<b>93±5</b>
cube-single-noisy-task1-v0	NOISY	<b>100±0</b>	<b>100±0</b>	<b>100±0</b>	1±1	93±6	<b>100±0</b>	<b>100±0</b>
cube-single-noisy-task2-v0	NOISY	<b>100±0</b>	<b>100±0</b>	<b>99±1</b>	7±6	92±9	<b>100±0</b>	<b>100±0</b>
cube-single-noisy-task3-v0	NOISY	<b>100±0</b>	<b>100±0</b>	<b>100±0</b>	85±9	53±25	<b>100±0</b>	<b>100±0</b>
cube-single-noisy-task4-v0	NOISY	<b>100±0</b>	<b>100±0</b>	<b>100±0</b>	0±0	91±4	<b>100±0</b>	<b>100±0</b>
cube-single-noisy-task5-v0	NOISY	<b>100±0</b>	<b>100±0</b>	<b>99±2</b>	0±0	91±5	<b>100±0</b>	<b>100±0</b>
cube-double-play-task1-v0	CLEAN	78±2	<b>96±2</b>	<b>97±1</b>	19±2	<b>95±1</b>	91±7	<b>97±1</b>
cube-double-play-task2-v0	CLEAN	52±3	83±11	83±1	15±4	63±4	83±7	<b>89±1</b>
cube-double-play-task3-v0	CLEAN	40±6	69±11	75±6	9±7	63±17	73±7	<b>81±7</b>
cube-double-play-task4-v0	CLEAN	7±3	21±1	20±4	1±1	17±5	22±4	<b>49±12</b>
cube-double-play-task5-v0	CLEAN	37±10	<b>81±10</b>	73±2	8±3	73±1	54±13	70±6
cube-double-noisy-task1-v0	NOISY	53±17	89±4	<b>97±2</b>	1±1	87±5	80±3	<b>92±2</b>
cube-double-noisy-task2-v0	NOISY	5±2	74±3	<b>80±4</b>	0±0	<b>79±11</b>	29±14	<b>81±10</b>
cube-double-noisy-task3-v0	NOISY	5±2	55±5	53±6	0±0	<b>60±5</b>	29±3	<b>61±12</b>
cube-double-noisy-task4-v0	NOISY	1±1	43±13	51±2	1±1	58±18	9±2	<b>65±5</b>
cube-double-noisy-task5-v0	NOISY	2±3	35±4	30±4	0±0	<b>40±11</b>	11±6	<b>38±13</b>
scene-play-task1-v0	CLEAN	<b>100±0</b>	<b>100±0</b>	<b>100±0</b>	<b>99±1</b>	<b>100±0</b>	<b>100±0</b>	<b>100±0</b>
scene-play-task2-v0	CLEAN	<b>100±0</b>	<b>100±0</b>	<b>99±1</b>	3±2	<b>100±0</b>	<b>100±0</b>	<b>100±0</b>
scene-play-task3-v0	CLEAN	93±5	93±5	94±3	80±4	<b>99±1</b>	93±8	<b>100±0</b>
scene-play-task4-v0	CLEAN	66±13	16±19	7±4	3±1	<b>100±0</b>	60±42	58±29
scene-play-task5-v0	CLEAN	<b>13±9</b>	0±0	0±0	0±0	0±0	0±0	0±0
scene-noisy-task1-v0	NOISY	86±12	<b>100±0</b>	<b>100±0</b>	<b>99±2</b>	<b>100±0</b>	71±17	<b>100±0</b>
scene-noisy-task2-v0	NOISY	70±8	<b>98±2</b>	<b>98±2</b>	0±0	<b>96±3</b>	<b>97±4</b>	<b>100±0</b>
scene-noisy-task3-v0	NOISY	0±0	25±6	39±16	44±4	<b>100±0</b>	45±27	80±27
scene-noisy-task4-v0	NOISY	0±0	8±11	1±1	1±1	37±39	0±0	<b>67±31</b>
scene-noisy-task5-v0	NOISY	<b>0±0</b>	<b>0±0</b>	<b>0±0</b>	<b>0±0</b>	<b>0±0</b>	<b>0±0</b>	<b>0±0</b>

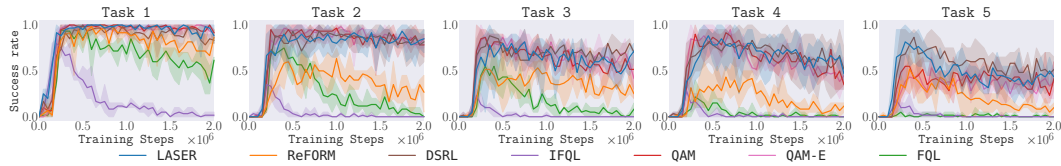


Figure 14: Training curves for the cube-double environment with the NOISY dataset.



Figure 15: Training curves for the scene environment with the CLEAN dataset.

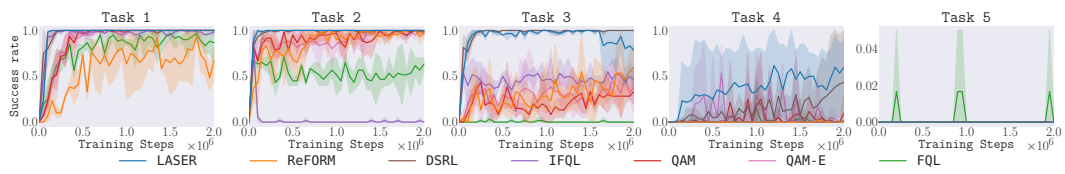


Figure 16: Training curves for the scene environment with the NOISY dataset.

## 800 **NeurIPS Paper Checklist**

### 801 **1. Claims**

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

804 Answer: [Yes]

805 Justification: All the claims have been justified in the paper.

806 Guidelines:

- 807 • The answer [N/A] means that the abstract and introduction do not include the claims made in  
808 the paper.
- 809 • The abstract and/or introduction should clearly state the claims made, including the contributions  
810 made in the paper and important assumptions and limitations. A [No] or [N/A] answer to this  
811 question will not be perceived well by the reviewers.
- 812 • The claims made should match theoretical and experimental results, and reflect how much the  
813 results can be expected to generalize to other settings.
- 814 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not  
815 attained by the paper.

### 816 **2. Limitations**

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

818 Answer: [Yes]

819 Justification: The limitations are discussed in Section 6.

820 Guidelines:

- 821 • The answer [N/A] means that the paper has no limitation while the answer [No] means that the  
822 paper has limitations, but those are not discussed in the paper.
- 823 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 824 • The paper should point out any strong assumptions and how robust the results are to violations of  
825 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,  
826 asymptotic approximations only holding locally). The authors should reflect on how these  
827 assumptions might be violated in practice and what the implications would be.
- 828 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested  
829 on a few datasets or with a few runs. In general, empirical results often depend on implicit  
830 assumptions, which should be articulated.
- 831 • The authors should reflect on the factors that influence the performance of the approach. For  
832 example, a facial recognition algorithm may perform poorly when image resolution is low or  
833 images are taken in low lighting. Or a speech-to-text system might not be used reliably to  
834 provide closed captions for online lectures because it fails to handle technical jargon.
- 835 • The authors should discuss the computational efficiency of the proposed algorithms and how  
836 they scale with dataset size.
- 837 • If applicable, the authors should discuss possible limitations of their approach to address  
838 problems of privacy and fairness.
- 839 • While the authors might fear that complete honesty about limitations might be used by reviewers  
840 as grounds for rejection, a worse outcome might be that reviewers discover limitations that  
841 aren’t acknowledged in the paper. The authors should use their best judgment and recognize  
842 that individual actions in favor of transparency play an important role in developing norms that  
843 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize  
844 honesty concerning limitations.

### 845 **3. Theory assumptions and proofs**

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

848 Answer: [Yes]

849 Justification: A full set of assumptions and complete and correct proof are provided in the  
850 appendix.

851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental result reproducibility

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

Answer: [Yes]

Justification: We have discussed the implementation details and provided hyperparameters in the appendix. We also include code with the submission.

Guidelines:

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

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We have included code in the supplementary material along with the commands.

Guidelines:

- 904 • The answer [N/A] means that paper does not include experiments requiring code.
- 905 • Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 906
- 907 • While we encourage the release of code and data, we understand that this might not be possible,
- 908 so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless
- 909 this is central to the contribution (e.g., for a new open-source benchmark).
- 910 • The instructions should contain the exact command and environment needed to run to reproduce
- 911 the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 912
- 913 • The authors should provide instructions on data access and preparation, including how to access
- 914 the raw data, preprocessed data, intermediate data, and generated data, etc.
- 915 • The authors should provide scripts to reproduce all experimental results for the new proposed
- 916 method and baselines. If only a subset of experiments are reproducible, they should state which
- 917 ones are omitted from the script and why.
- 918 • At submission time, to preserve anonymity, the authors should release anonymized versions (if
- 919 applicable).
- 920 • Providing as much information as possible in supplemental material (appended to the paper) is
- 921 recommended, but including URLs to data and code is permitted.

## 922 6. Experimental setting/details

923 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,

924 how they were chosen, type of optimizer) necessary to understand the results?

925 Answer: [Yes]

926 Justification: All details are provided in the appendix.

927 Guidelines:

- 928 • The answer [N/A] means that the paper does not include experiments.
- 929 • The experimental setting should be presented in the core of the paper to a level of detail that is
- 930 necessary to appreciate the results and make sense of them.
- 931 • The full details can be provided either with the code, in appendix, or as supplemental material.

## 932 7. Experiment statistical significance

933 Question: Does the paper report error bars suitably and correctly defined or other appropriate

934 information about the statistical significance of the experiments?

935 Answer: [Yes]

936 Justification: We report error bars in all figures and tables that contain experimental results.

937 Guidelines:

- 938 • The answer [N/A] means that the paper does not include experiments.
- 939 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
- 940 intervals, or statistical significance tests, at least for the experiments that support the main claims
- 941 of the paper.
- 942 • The factors of variability that the error bars are capturing should be clearly stated (for example,
- 943 train/test split, initialization, random drawing of some parameter, or overall run with given
- 944 experimental conditions).
- 945 • The method for calculating the error bars should be explained (closed form formula, call to a
- 946 library function, bootstrap, etc.)
- 947 • The assumptions made should be given (e.g., Normally distributed errors).
- 948 • It should be clear whether the error bar is the standard deviation or the standard error of the
- 949 mean.
- 950 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably
- 951 report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of
- 952 errors is not verified.
- 953 • For asymmetric distributions, the authors should be careful not to show in tables or figures
- 954 symmetric error bars that would yield results that are out of range (e.g., negative error rates).

- 955 • If error bars are reported in tables or plots, the authors should explain in the text how they were  
956 calculated and reference the corresponding figures or tables in the text.

957 **8. Experiments compute resources**

958 Question: For each experiment, does the paper provide sufficient information on the computer  
959 resources (type of compute workers, memory, time of execution) needed to reproduce the experi-  
960 ments?

961 Answer: [Yes]

962 Justification: We have included this information in the appendix.

963 Guidelines:

- 964 • The answer [N/A] means that the paper does not include experiments.  
965 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud  
966 provider, including relevant memory and storage.  
967 • The paper should provide the amount of compute required for each of the individual experimental  
968 runs as well as estimate the total compute.  
969 • The paper should disclose whether the full research project required more compute than the  
970 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it  
971 into the paper).

972 **9. Code of ethics**

973 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS  
974 Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

975 Answer: [Yes]

976 Justification: The research conducted in the paper conforms, in every respect, with the NeurIPS  
977 Code of Ethics.

978 Guidelines:

- 979 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.  
980 • If the authors answer [No], they should explain the special circumstances that require a deviation  
981 from the Code of Ethics.  
982 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due  
983 to laws or regulations in their jurisdiction).

984 **10. Broader impacts**

985 Question: Does the paper discuss both potential positive societal impacts and negative societal  
986 impacts of the work performed?

987 Answer: [Yes]

988 Justification: The work does not have negative societal impacts.

989 Guidelines:

- 990 • The answer [N/A] means that there is no societal impact of the work performed.  
991 • If the authors answer [N/A] or [No], they should explain why their work has no societal impact  
992 or why the paper does not address societal impact.  
993 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,  
994 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-  
995 ment of technologies that could make decisions that unfairly impact specific groups), privacy  
996 considerations, and security considerations.  
997 • The conference expects that many papers will be foundational research and not tied to par-  
998 ticular applications, let alone deployments. However, if there is a direct path to any negative  
999 applications, the authors should point it out. For example, it is legitimate to point out that  
1000 an improvement in the quality of generative models could be used to generate Deepfakes for  
1001 disinformation. On the other hand, it is not needed to point out that a generic algorithm for  
1002 optimizing neural networks could enable people to train models that generate Deepfakes faster.  
1003 • The authors should consider possible harms that could arise when the technology is being used  
1004 as intended and functioning correctly, harms that could arise when the technology is being used  
1005 as intended but gives incorrect results, and harms following from (intentional or unintentional)  
1006 misuse of the technology.

- 1007 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies  
1008 (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for  
1009 monitoring misuse, mechanisms to monitor how a system learns from feedback over time,  
1010 improving the efficiency and accessibility of ML).

## 1011 11. Safeguards

1012 Question: Does the paper describe safeguards that have been put in place for responsible release  
1013 of data or models that have a high risk for misuse (e.g., pre-trained language models, image  
1014 generators, or scraped datasets)?

1015 Answer: [N/A]

1016 Justification: The paper does not contain such data or models.

1017 Guidelines:

- 1018 • The answer [N/A] means that the paper poses no such risks.
- 1019 • Released models that have a high risk for misuse or dual-use should be released with necessary  
1020 safeguards to allow for controlled use of the model, for example by requiring that users adhere  
1021 to usage guidelines or restrictions to access the model or implementing safety filters.
- 1022 • Datasets that have been scraped from the Internet could pose safety risks. The authors should  
1023 describe how they avoided releasing unsafe images.
- 1024 • We recognize that providing effective safeguards is challenging, and many papers do not require  
1025 this, but we encourage authors to take this into account and make a best faith effort.

## 1026 12. Licenses for existing assets

1027 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the  
1028 paper, properly credited and are the license and terms of use explicitly mentioned and properly  
1029 respected?

1030 Answer: [Yes]

1031 Justification: We have cited the original papers.

1032 Guidelines:

- 1033 • The answer [N/A] means that the paper does not use existing assets.
- 1034 • The authors should cite the original paper that produced the code package or dataset.
- 1035 • The authors should state which version of the asset is used and, if possible, include a URL.
- 1036 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1037 • For scraped data from a particular source (e.g., website), the copyright and terms of service of  
1038 that source should be provided.
- 1039 • If assets are released, the license, copyright information, and terms of use in the package should  
1040 be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for  
1041 some datasets. Their licensing guide can help determine the license of a dataset.
- 1042 • For existing datasets that are re-packaged, both the original license and the license of the derived  
1043 asset (if it has changed) should be provided.
- 1044 • If this information is not available online, the authors are encouraged to reach out to the asset's  
1045 creators.

## 1046 13. New assets

1047 Question: Are new assets introduced in the paper well documented and is the documentation  
1048 provided alongside the assets?

1049 Answer: [Yes]

1050 Justification: We provide the commands along with the code.

1051 Guidelines:

- 1052 • The answer [N/A] means that the paper does not release new assets.
- 1053 • Researchers should communicate the details of the dataset/code/model as part of their sub-  
1054 missions via structured templates. This includes details about training, license, limitations,  
1055 etc.
- 1056 • The paper should discuss whether and how consent was obtained from people whose asset is  
1057 used.

- 1058 • At submission time, remember to anonymize your assets (if applicable). You can either create  
1059 an anonymized URL or include an anonymized zip file.

1060 **14. Crowdsourcing and research with human subjects**

1061 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1062 include the full text of instructions given to participants and screenshots, if applicable, as well as  
1063 details about compensation (if any)?

1064 Answer: [N/A]

1065 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1066 Guidelines:

- 1067 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with  
1068 human subjects.
- 1069 • Including this information in the supplemental material is fine, but if the main contribution of  
1070 the paper involves human subjects, then as much detail as possible should be included in the  
1071 main paper.
- 1072 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other  
1073 labor should be paid at least the minimum wage in the country of the data collector.

1074 **15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

1075 Question: Does the paper describe potential risks incurred by study participants, whether such  
1076 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals  
1077 (or an equivalent approval/review based on the requirements of your country or institution) were  
1078 obtained?

1079 Answer: [N/A]

1080 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1081 Guidelines:

- 1082 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with  
1083 human subjects.
- 1084 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be  
1085 required for any human subjects research. If you obtained IRB approval, you should clearly  
1086 state this in the paper.
- 1087 • We recognize that the procedures for this may vary significantly between institutions and  
1088 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for  
1089 their institution.
- 1090 • For initial submissions, do not include any information that would break anonymity (if applica-  
1091 ble), such as the institution conducting the review.

1092 **16. Declaration of LLM usage**

1093 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-  
1094 standard component of the core methods in this research? Note that if the LLM is used only for  
1095 writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor,  
1096 or originality of the research, declaration is not required.

1097 Answer: [N/A]

1098 Justification: The core method development in this research does not involve LLMs as any  
1099 important, original, or non-standard components.

1100 Guidelines:

- 1101 • The answer [N/A] means that the core method development in this research does not involve  
1102 LLMs as any important, original, or non-standard components.
- 1103 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not be  
1104 described.