

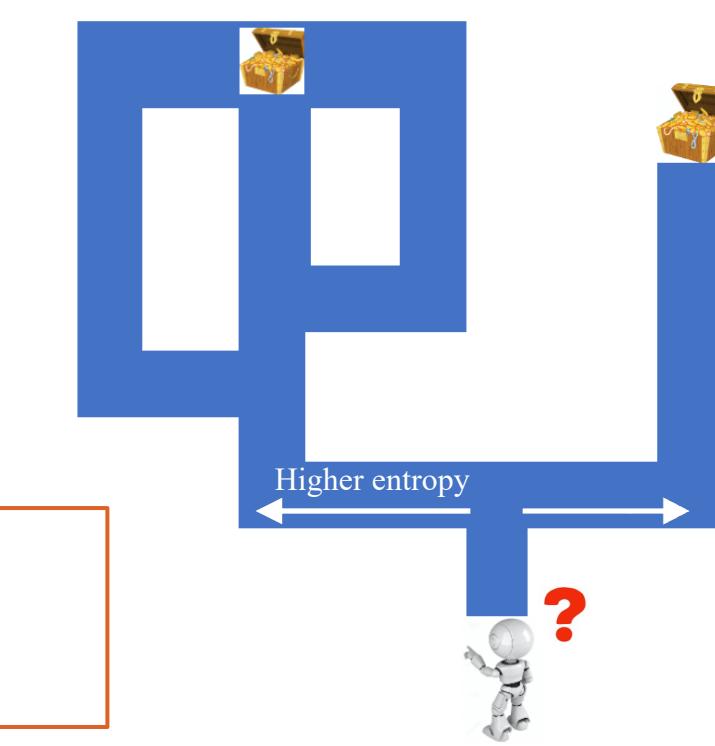
Motivation: A more optimal solution to MaxEntr RL

MaxEntr RL searches for the policy that maximizes the expected **future entropy** jointly with the expected **future reward**

$$\pi^* = \operatorname{argmax}_{\pi \in \Pi} \sum_t \gamma^t \mathbb{E}_{(s_t, a_t) \sim \rho_\pi} [r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot | s_t))]$$

Why MaxEntr RL?

- Provably robust
- Better exploration
- Training stability
- Sample efficiency



Open Challenge
How can we estimate the entropy of arbitrarily complex policies?

Landmark MaxEntr RL Algorithms

- Soft Q-learning (SQL) [Haarnoja *et al.*, ICML 2017]
 - ✓ Implicit entropy computation for arbitrary policies
 - ✗ Not scalable (uses importance sampling)
- Soft Actor Critic (SAC) [Haarnoja *et al.*, ICML 2018]
 - ✓ Tractable entropy estimate (policy=Gaussian)
 - ✗ Limited to unimodal policies
- Soft Actor Critic with Normalizing Flows (SAC-NF) [Mazouze *et al.*, 2020]
 - ✓ Tractable entropy estimate (policy=Normalizing Flow)
 - ✗ Hard to train (collapses to local minima)

Our Approach: Stein Soft Actor Critic (S²AC)

- Policy = Stein Variational Gradient Descent (SVGD) Sampler [Liu *et al.*, 2017]
 - ✓ **Closed form expression of the entropy of arbitrary complex distribution**
 - ✓ Computationally efficient (only vector dot products and first order derivatives)
 - ✓ Parameter efficient (same num. parameters as SAC)

Background

Change of Variable Formula (CVF):

$$F: Z \rightarrow X \text{ is invertible} \Rightarrow p_X(x) = p_Z(z) \left| \det \frac{\partial F(z)}{\partial z} \right|^{-1}$$

S²AC: Stein Soft Actor Critic

Optimal policy:

$$\pi^*(a|s) = \frac{\exp(Q_\phi(s, a)/\alpha)}{Z} \quad \boxed{\text{Untractable}}$$

Critic: Learn the **Q-values** by optimizing the Bellman loss

$$\phi^* = \arg \min_{\phi} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi_\theta}} [(Q_\phi(s_t, a_t) - \hat{y})^2]$$

$$\hat{y} = r_t(s_t, a_t) + \gamma \mathbb{E}_{(s_{t+1}, a_{t+1}) \sim \rho_{\pi_\theta}} [Q_\phi(s_{t+1}, a_{t+1}) + \alpha \mathcal{H}(\pi_\theta(\cdot | s_{t+1}))]$$

Actor: Learn a **policy (SVGD sampler)** via variational inference

$$\theta^* = \operatorname{argmin}_\theta D_{KL}\left(\pi_\theta(s|a) \parallel \frac{\exp(Q_\phi(s, a)/\alpha)}{Z}\right)$$

SVGD Update Rule (step):

$$a_i^{l+1} = a_i^l + \epsilon \mathbb{E}_{a_j^l \sim q_\theta^l} \left[k(a_i^l, a_j^l) \nabla_{a_j^l} \log \pi^*(a_j^l | s) + \nabla_{a_j^l} k(a_i^l, a_j^l) \right]$$

Weighted average of particles gradients Deterministic repulsion

SVGD distribution q_θ^0 CVF q_θ^1 CVF q_θ^2 CVF $q_\theta^L = \pi_\theta$

$k(\cdot, \cdot)$: kernel l : SVGD iteration

A Closed Form Estimate of The SVGD Induced Distribution:

The closed-form estimate of $\log q^L(a^L | s)$ for the SVGD based sampler with an RBF kernel is:

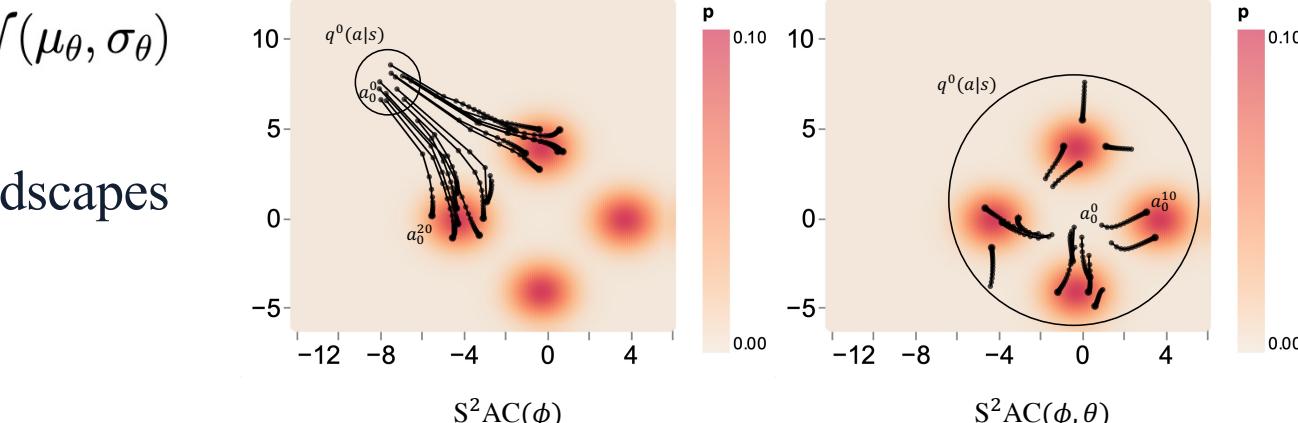
$$\log q^L(a^L | s) \approx \log q^0(a^0 | s) + \frac{\epsilon}{m\sigma^2} \sum_{l=0}^{L-1} \sum_{j=1, a^l \neq a_j^l}^m k(a_j^l, a^l) \left((a^l - a_j^l)^\top \nabla_{a_j^l} Q(s, a_j^l) + \frac{\alpha}{\sigma^2} \|a^l - a_j^l\|^2 - d\alpha \right)$$

neighborhood Curvature dispersion

Parametrized Initial Distribution $q_\theta^0 = \mathcal{N}(\mu_\theta, \sigma_\theta)$

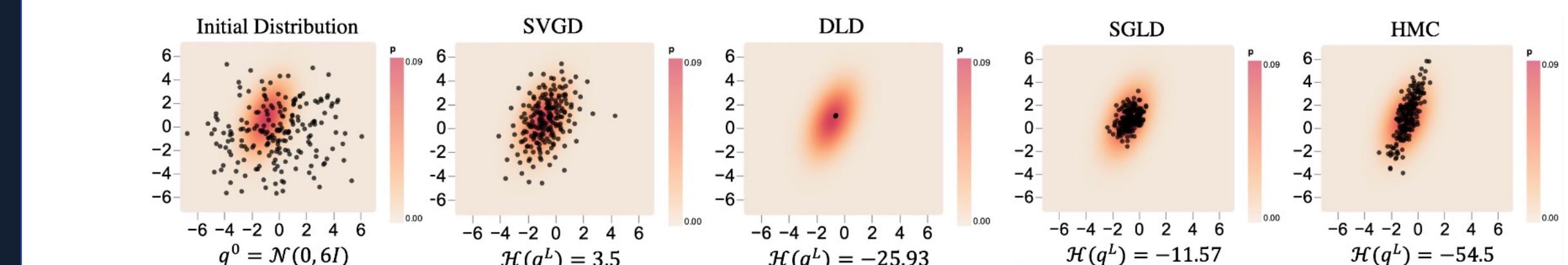
- Faster convergence
- Robustness to non-smooth deepnet landscapes by spatially constraining the actions:

$$(\mu - 3\sigma_\theta \leq a_\theta \leq \mu + 3\sigma_\theta)$$



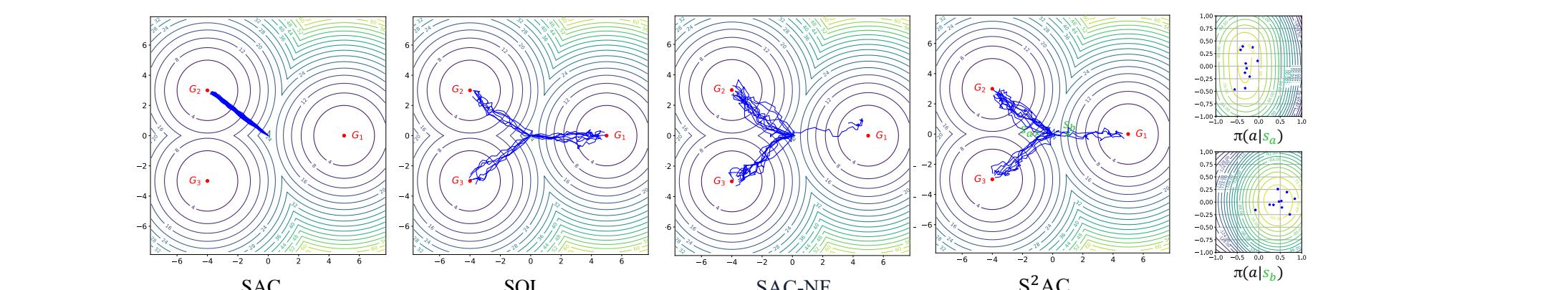
Results

Entropy Evaluation ($\mathcal{H}(\pi^*) = 3.41$)

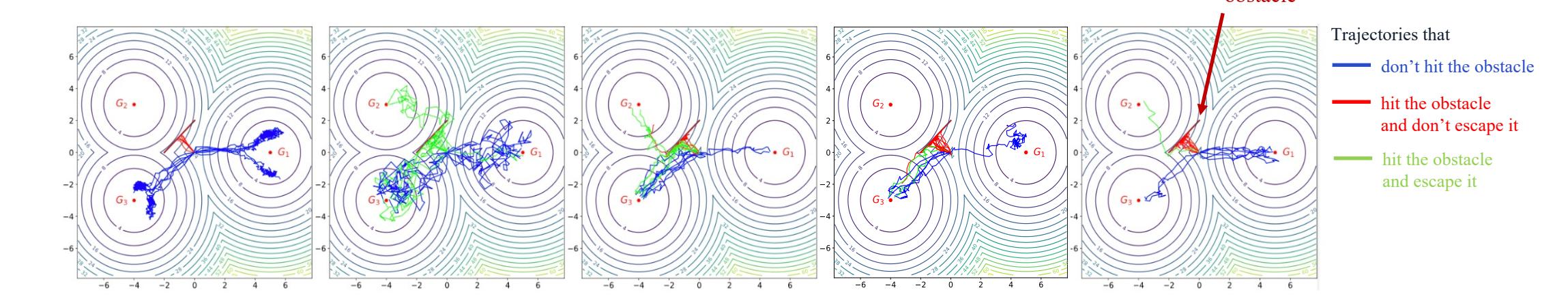


Multigoal Environment

- Multi-Modality: S²AC maximizes the expected future entropy and recovers all three modes



- Robustness: S²AC escapes obstacles at test time more frequently than other baselines



MuJoCo Environment

- Performance: S²AC beats the baselines on 4 out of 5 environments
- Run-time: Amortized S²AC has comparable run-time to SAC w/o trading-off performance

