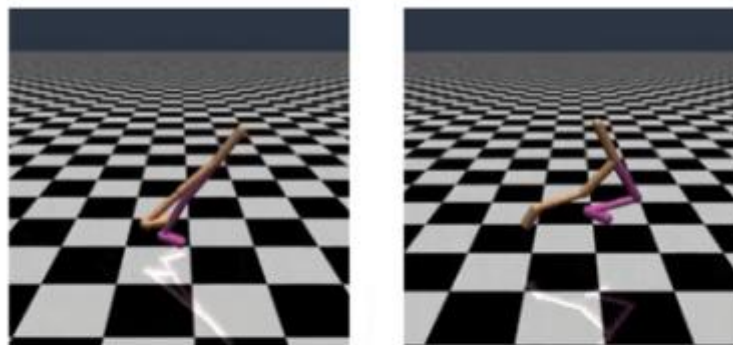


---

## Reinforcement Learning with Deep Energy-Based Policies

---

Tuomas Haarnoja<sup>\*1</sup> Haoran Tang<sup>\*2</sup> Pieter Abbeel<sup>1,3,4</sup> Sergey Levine<sup>1</sup>

 $\pi_{\theta_1}$  $\pi_{\theta_2}$ 

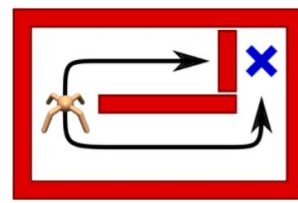
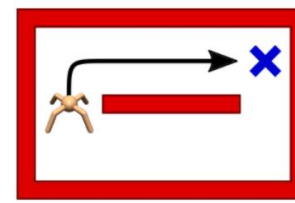
$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^T r_t \right]$$

Exploration Problem.

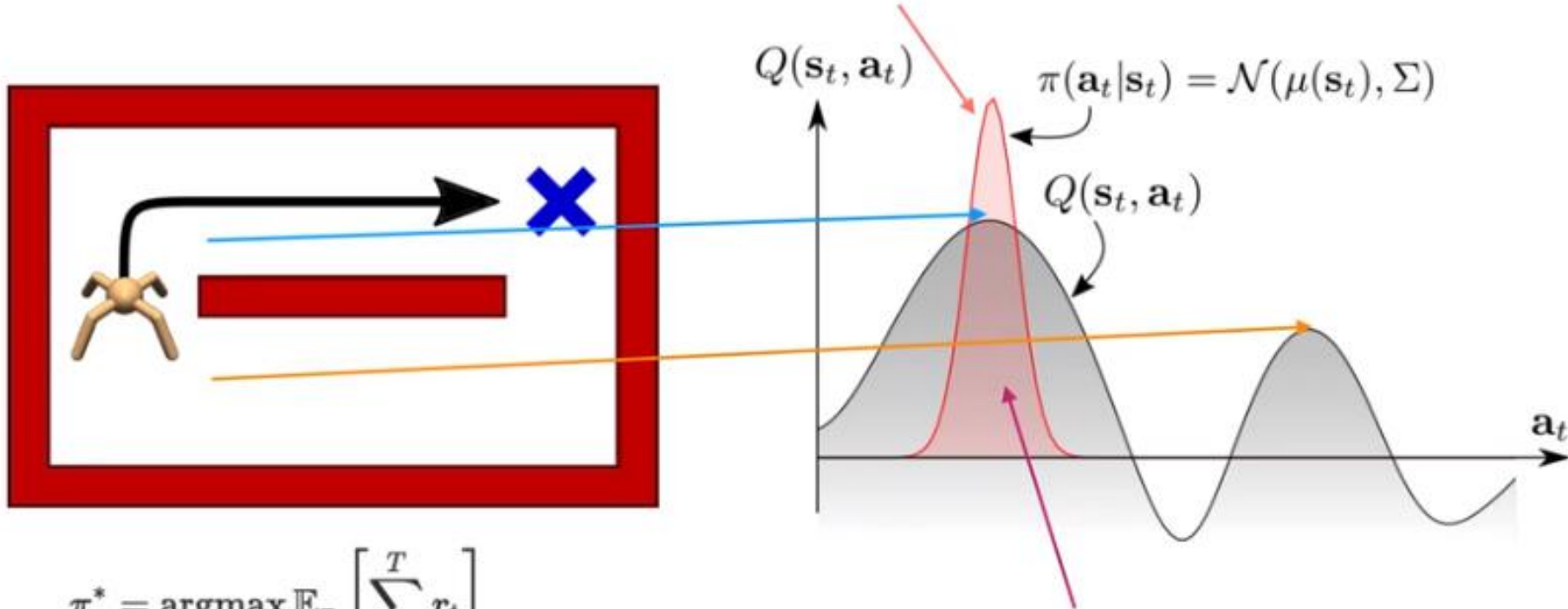
[Motivation]

Standard RL aims to master **a single way** to solve a given task.

# Problem of Standard RL



*Unimodal Policy: only walk up*



$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^T r_t \right]$$

Gaussian policy

$$\pi(a_t|s_t) = \mathcal{N}(\mu(s_t), \Sigma)$$

- Standard RL objective finds a unimodal policy
- center at the maximal Q-value.
- add noise for exploration.

# Energy-based Policy

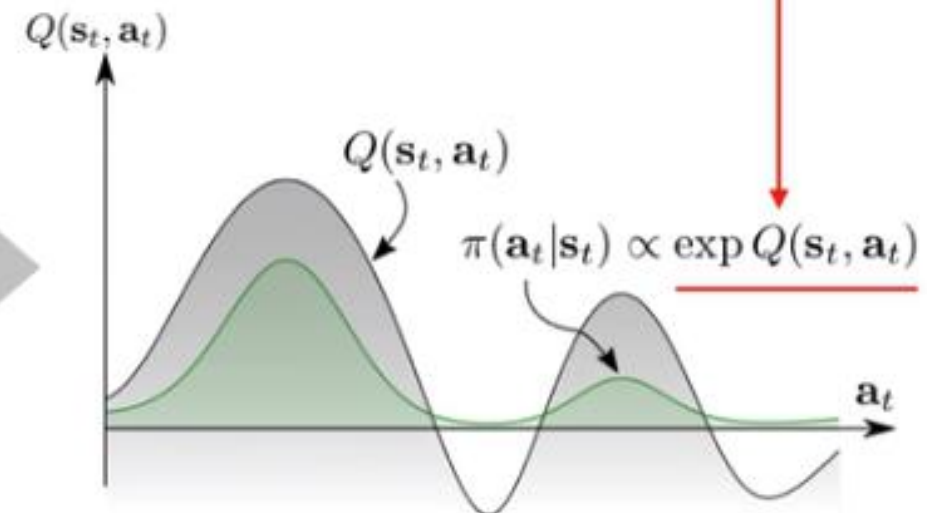
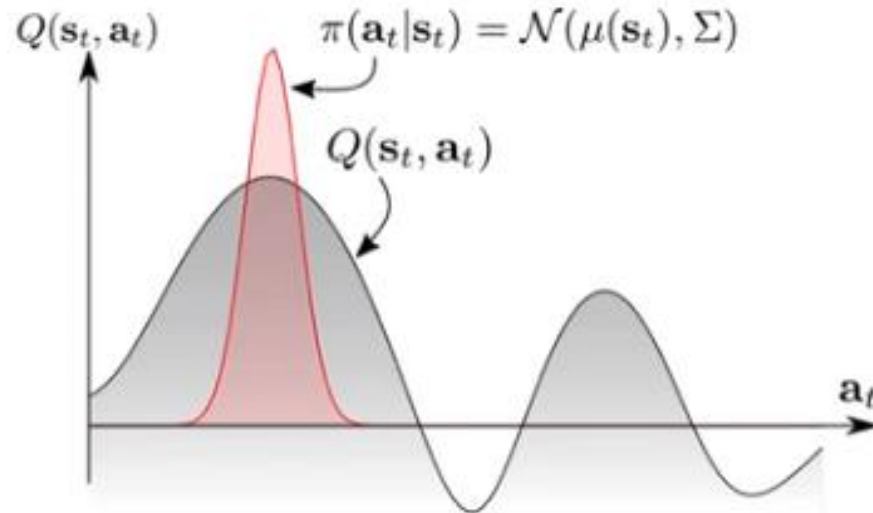
[Intuition] Exploring all promising states while prioritizing the best one.

Define the policy with exponentiated Q-values ==> Boltzmann Distribution

$$\pi(\mathbf{a}_t | \mathbf{s}_t) \propto \exp(-\mathcal{E}(\mathbf{s}_t, \mathbf{a}_t))$$

$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^T r_t \right]$$

$$\mathcal{E}(\mathbf{s}_t, \mathbf{a}_t) = -\frac{1}{\alpha} Q(\mathbf{s}_t, \mathbf{a}_t)$$



# Entropy Regularization

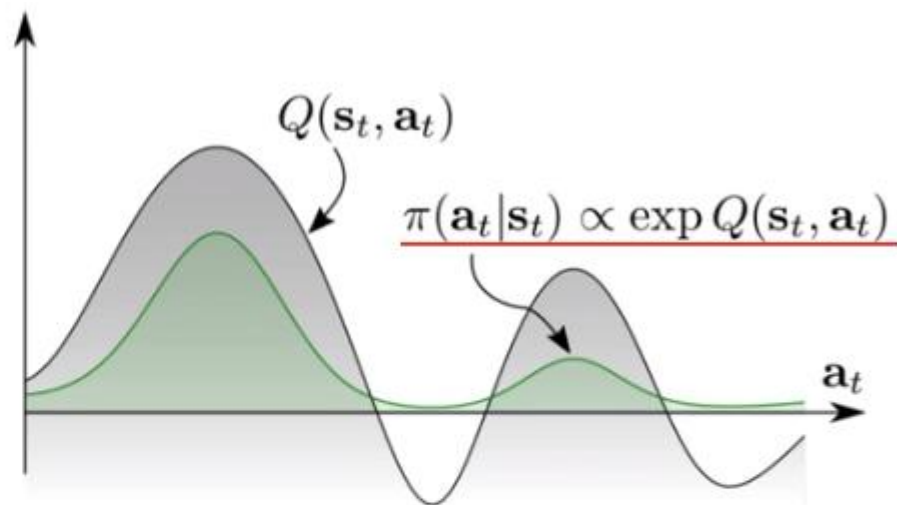
Boltzmann distribution policy:

$$\tilde{\pi}(a_t|s_t) = \frac{\exp(Q^\pi(s_t, a_t))}{\int \exp(Q^\pi(s_t, a)) da}$$

At state  $s_t$ , the KL divergence:

$$\begin{aligned} D_{KL} [\pi(\cdot|s_t) \|\tilde{\pi}(\cdot|s_t)] &= \mathbb{E}_{a_t \sim \pi} [\log \pi(\cdot|s_t) - Q^\pi(s_t, a_t) + \log \int \exp(Q^\pi(s_t, a)) da] \\ &= -\mathcal{H}(\pi(\cdot|s_t)) - \mathbb{E}_{a_t \sim \pi} [Q^\pi(s_t, a_t)] + \log \int \exp(Q^\pi(s_t, a)) da \end{aligned}$$

$$\underbrace{\mathcal{H}(\pi(\cdot|s_t)) + \mathbb{E}_{a_t \sim \pi} [Q^\pi(s_t, a_t)]}_{\text{objective}} = \underbrace{-D_{KL} [\pi(\cdot|s_t) \|\tilde{\pi}(\cdot|s_t)]}_{\text{entropy regularization}} + \underbrace{\log \int \exp(Q^\pi(s_t, a)) da}_{\text{partition function}}$$



# Soft Q-Value

Soft Q-value ( $\alpha = 1$ ) is the expectation of the discounted sum of rewards and entropy:

$$Q^\pi(\mathbf{s}, \mathbf{a}) \triangleq r_0 + \mathbb{E}_{\tau \sim \pi, \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a}} \left[ \sum_{t=1}^{\infty} \gamma^t (r_t + \underbrace{\mathcal{H}(\pi(\cdot | \mathbf{s}_t))}_{\text{entropy regularized reward}}) \right]$$

$$\begin{aligned} Q^\pi(\mathbf{s}, \mathbf{a}) &= r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p} [\mathcal{H}(\pi(\cdot | \mathbf{s}')) + \mathbb{E}_{\mathbf{a}' \sim \pi(\cdot | \mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')] ] \\ &= r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p} [V^\pi(\mathbf{s}')] . \end{aligned}$$

soft-max

$$V^\pi(\mathbf{s}) \triangleq \log \int \exp(Q^\pi(\mathbf{s}, \mathbf{a})) \, d\mathbf{a} \quad \approx \quad V^\pi(s) = \max_a Q^\pi(s, a)$$

hard-max

$$\pi(a_t | s_t) = \frac{\exp(Q^\pi(s_t, a_t))}{\int \exp(Q^\pi(s_t, a)) da} = \exp(Q^\pi(\mathbf{s}, \mathbf{a}) - V^\pi(\mathbf{s}))$$

soft Bellman operator

$$\mathcal{T}Q(\mathbf{s}, \mathbf{a}) \triangleq r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p} \left[ \log \int \exp Q(\mathbf{s}', \mathbf{a}') \, d\mathbf{a}' \right]$$

Bellman operator

$$\mathcal{T}Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p} \left[ \max_{a'} Q(s', a') \right]$$

# Soft Q-Learning

Soft Bellman Backup Operator  $\mathcal{T}$  Contraction

**Theorem 1.** Soft Q-iteration. Let  $Q$  and  $V$  be bounded and  $\int_{\mathcal{A}} \exp(\frac{1}{\alpha} Q(s_t, \mathbf{a}')) d\mathbf{a}' < \infty$ ,  $\forall s_t \in S$ , and assume that  $Q^* < \infty$  exists. Then the fixed-point iteration

$$Q(s_t, \mathbf{a}_t) \leftarrow r(s_t, \mathbf{a}_t) + \gamma \mathbb{E}_{s_{t+1} \sim p} [V(s_{t+1})], \quad \forall s_t, \mathbf{a}_t \quad (3.6)$$

$$V(s_t) \leftarrow \alpha \log \left( \int_{\mathcal{A}} \exp \left( \frac{1}{\alpha} Q(s_t, \mathbf{a}') \right) d\mathbf{a}' \right), \quad \forall s_t \quad (3.7)$$

converges to  $Q^*$  and  $V^*$ , respectively.

Boltzmann  
Distribution

Soft max

Entropy regularized  
reward

**Optimal Soft Q-function:**  $Q_{\text{soft}}^*(s_t, \mathbf{a}_t) = r_t + \mathbb{E}_{(s_{t+1}, \dots) \sim \rho_{\pi}} \left[ \sum_{l=1}^{\infty} \gamma^l (r_{t+l} + \alpha \mathcal{H}(\pi_{\text{MaxEnt}}^*(\cdot | s_{t+l}))) \right]$

**Optimal Soft value function:**  $V^*(s_t) = \alpha \log \int_{\mathcal{A}} \exp \left( \frac{1}{\alpha} Q^*(s_t, \mathbf{a}') \right) d\mathbf{a}' \longleftrightarrow V^{\pi}(s) = \max_a Q^{\pi}(s, a)$

**Optimal policy:**  $\pi_{\text{MaxEnt}}^*(\mathbf{a}_t | s_t) = \exp \left( \frac{1}{\alpha} (Q_{\text{soft}}^*(s_t, \mathbf{a}_t) - V_{\text{soft}}^*(s_t)) \right) \longleftrightarrow \pi(a | s) = \max_a Q(s, a)$



# Problems of Soft Q-learning

$$\pi(a_t|s_t) = \frac{\exp(Q^\pi(s_t, a_t))}{\int \exp(Q^\pi(s_t, a)) da}$$

Hard to sample from  $\pi$

Soft Bellman Backup Operator

$$\begin{aligned} Q_{\text{soft}}(s_t, a_t) &\leftarrow r_t + \gamma \mathbb{E}_{s_{t+1} \sim p_{\pi}} [V_{\text{soft}}(s_{t+1})], \forall s_t, a_t \\ V_{\text{soft}}(s_t) &\leftarrow \alpha \log \int_{\mathcal{A}} \exp\left(\frac{1}{\alpha} Q_{\text{soft}}(s_t, a')\right) da', \forall s_t \end{aligned}$$

Integral over infinite states and actions

$$\underbrace{V^*(s_t) = \alpha \log \int_{\mathcal{A}} \exp\left(\frac{1}{\alpha} Q^*(s_t, a')\right) da'}_{\text{intractable}} \xrightarrow{\text{IS}} V_\theta(s_t) = \alpha \log \mathbb{E}_{q_{a'}} \left[ \frac{\exp\left(\frac{1}{\alpha} Q_\theta(s_t, a')\right)}{q_{a'}(a')} \right]$$

arbitrary distribution

Minimize soft Bellman error  $|\mathcal{T}Q - Q| \Rightarrow \text{sample } (s_t, a_t, s_{t+1}) \text{ from ER and SGD}$

$$J_Q(\theta) = \underbrace{\mathbb{E}_{s_t \sim q_s, a_t \sim q_a}}_{\text{arbitrary sampling distribution}} \left[ \frac{1}{2} \left( Q_\theta(s_t, a_t) - \underbrace{\left( r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim p} [V_{\bar{\theta}}(s_{t+1})] \right)}_{\text{target value network}} \right)^2 \right]$$

# Approximate Sampling and Stein Variational Gradient Descent (SVGD)

$$\pi(a_t|s_t) = \frac{\exp(Q^\pi(s_t, a_t))}{\int \exp(Q^\pi(s_t, a)) da}$$

A stochastic NN  $\mathbf{a}_t = f^\phi(\xi; \mathbf{s}_t)$  maps noise  $\xi$  to action samples.

Stochastic Sampling  
Network

**Objective**  $J_\pi(\phi; \mathbf{s}_t) = D_{\text{KL}} \left( \pi_\phi(\cdot | \mathbf{s}_t) \parallel \exp \left( \frac{1}{\alpha} (Q_\theta(\mathbf{s}_t, \cdot) - V_\theta(\mathbf{s}_t)) \right) \right)$

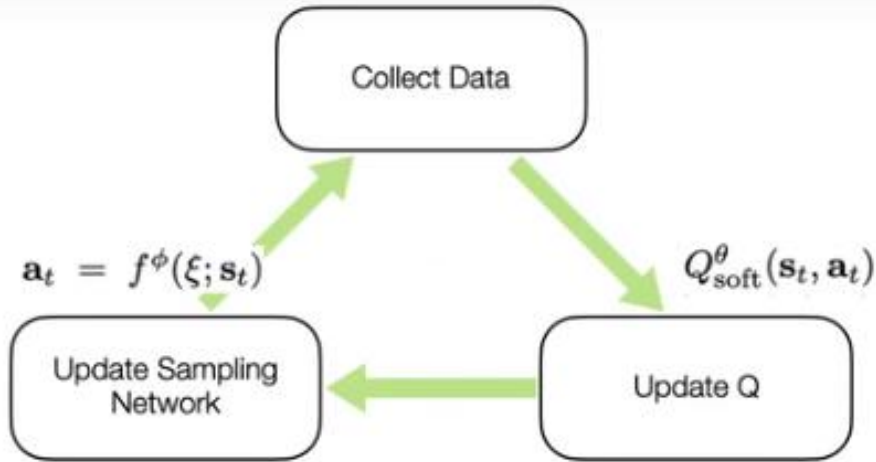
Suppose we have a set of independent, state conditioned action samples  $\{\mathbf{a}_t^{(i)} = f_\phi(\xi^{(i)}; \mathbf{s}_t)\}$ , given by the sampling network, and we “perturb” them in some directions  $\Delta f_\phi(\xi^{(i)}; \mathbf{s}_t)$ , then the KL divergence induced by the samples can be reduced. Stein variational gradient descent (Liu & Wang, 2016) provides the **most greedy directions** as a functional

$$\Delta f_\phi(\cdot; \mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} \left[ \kappa(\mathbf{a}_t, f_\phi(\cdot; \mathbf{s}_t)) \nabla_{\mathbf{a}'} Q_\theta(\mathbf{s}_t, \mathbf{a}')|_{\mathbf{a}'=\mathbf{a}_t} + \alpha \nabla_{\mathbf{a}'} \kappa(\mathbf{a}', f_\phi(\cdot; \mathbf{s}_t))|_{\mathbf{a}'=\mathbf{a}_t} \right], \quad (3.11)$$

$$\frac{\partial J_\pi}{\partial \mathbf{a}_t} \propto \Delta f^\phi \longrightarrow \frac{\partial J_\pi(\phi; \mathbf{s}_t)}{\partial \phi} \propto \mathbb{E}_\xi \left[ \Delta f^\phi(\xi; \mathbf{s}_t) \frac{\partial f^\phi(\xi; \mathbf{s}_t)}{\partial \phi} \right]$$

Wang, D. and Liu, Q. Learning to draw samples: With application to amortized mle for generative adversarial learning. *arXiv preprint arXiv:1611.01722*, 2016.





$$V_{\theta}(s_t) = \alpha \log \mathbb{E}_{q_{a'}} \left[ \frac{\exp \left( \frac{1}{\alpha} Q_{\theta}(s_t, a') \right)}{q_{a'}(a')} \right]$$

$$J_Q(\theta) = \mathbb{E}_{s_t \sim q_s, a_t \sim q_a} \left[ \frac{1}{2} \left( Q_{\theta}(s_t, a_t) - \left( r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim p} [V_{\hat{\theta}}(s_{t+1})] \right) \right)^2 \right]$$

### Stein Variational Gradient Descent

$$\Delta f^{\phi}(\cdot; s_t) = \mathbb{E}_{a_t \sim \pi^{\phi}} \left[ \kappa(a_t, f^{\phi}(\cdot; s_t)) \nabla_{a'} Q_{\text{soft}}^{\theta}(s_t, a') \Big|_{a'=a_t} + \alpha \nabla_{a'} \kappa(a', f^{\phi}(\cdot; s_t)) \Big|_{a'=a_t} \right], \quad (13)$$

$$\frac{\partial J_{\pi}(\phi; s_t)}{\partial \phi} \propto \mathbb{E}_{\xi} \left[ \Delta f^{\phi}(\xi; s_t) \frac{\partial f^{\phi}(\xi; s_t)}{\partial \phi} \right]$$

### Algorithm 1 Soft Q-learning

$\theta, \phi \sim$  some initialization distributions.

Assign target parameters:  $\bar{\theta} \leftarrow \theta, \bar{\phi} \leftarrow \phi$ .

$\mathcal{D} \leftarrow$  empty replay memory.

**for each epoch do**

**for each  $t$  do**

**Collect experience**

Sample an action for  $s_t$  using  $f^{\phi}$ :

$a_t \leftarrow f^{\phi}(\xi; s_t)$  where  $\xi \sim \mathcal{N}(\mathbf{0}, I)$ .

Sample next state from the environment:

$s_{t+1} \sim p_s(s_{t+1} | s_t, a_t)$ .

Save the new experience in the replay memory:

$\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, a_t, r(s_t, a_t), s_{t+1})\}$ .

**Sample a minibatch from the replay memory**

$\{(s_t^{(i)}, a_t^{(i)}, r_t^{(i)}, s_{t+1}^{(i)})\}_{i=0}^N \sim \mathcal{D}$ .

**Update the soft Q-function parameters**

Sample  $\{a^{(i,j)}\}_{j=0}^M \sim q_{a'}$  for each  $s_{t+1}^{(i)}$ .

Compute empirical soft values  $\hat{V}_{\text{soft}}^{\bar{\theta}}(s_{t+1}^{(i)})$  in (10).

Compute empirical gradient  $\hat{\nabla}_{\theta} J_Q$  of (11).

Update  $\theta$  according to  $\hat{\nabla}_{\theta} J_Q$  using ADAM.

**Update policy**

Sample  $\{\xi^{(i,j)}\}_{j=0}^M \sim \mathcal{N}(\mathbf{0}, I)$  for each  $s_t^{(i)}$ .

Compute actions  $a_t^{(i,j)} = f^{\phi}(\xi^{(i,j)}, s_t^{(i)})$ .

Compute  $\Delta f^{\phi}$  using empirical estimate of (13).

Compute empirical estimate of (14):  $\hat{\nabla}_{\phi} J_{\pi}$ .

Update  $\phi$  according to  $\hat{\nabla}_{\phi} J_{\pi}$  using ADAM.

**end for**

**if epoch mod update\_interval = 0 then**

Update target parameters:  $\bar{\theta} \leftarrow \theta, \bar{\phi} \leftarrow \phi$ .

**end if**

**end for**

# Experiment

- (1) Can SQL capture a multi-modal policy distribution?
- (2) Can SQL help exploration?
- (3) Can SQL serve as a good initialization for fine-tuning on different tasks.

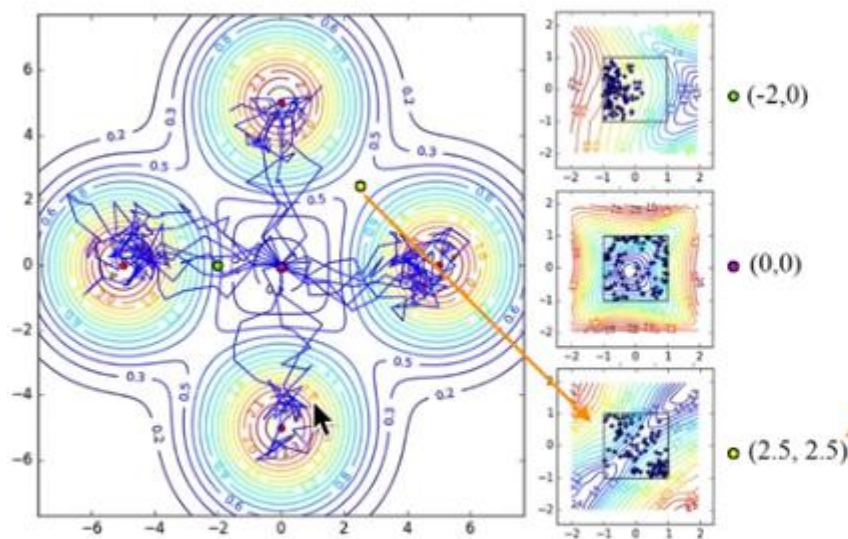
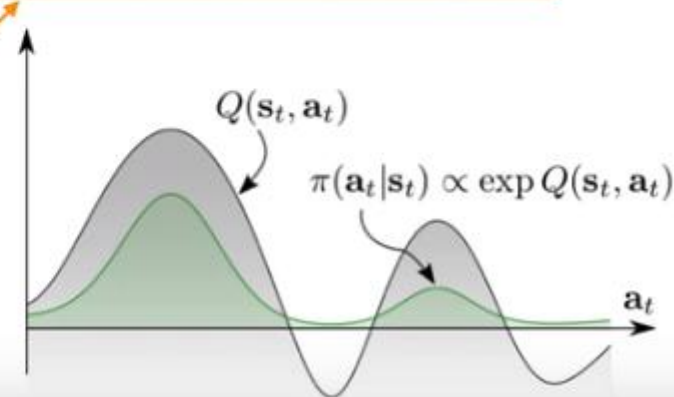
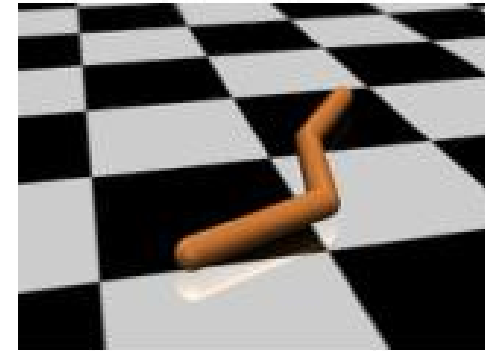


Figure 1. Illustration of the 2D multi-goal environment. Left: trajectories from a policy learned with our method (solid blue lines). The  $x$  and  $y$  axes correspond to 2D positions (states). The agent is initialized at the origin. The goals are depicted as red dots, and the level curves show the reward. Right: Q-values at three selected states, depicted by level curves (red: high values, blue: low values). The  $x$  and  $y$  axes correspond to 2D velocity (actions) bounded between -1 and 1. Actions sampled from the policy are shown as blue stars. Note that, in regions (e.g.  $(2.5, 2.5)$ ) between the goals, the method chooses multimodal actions.

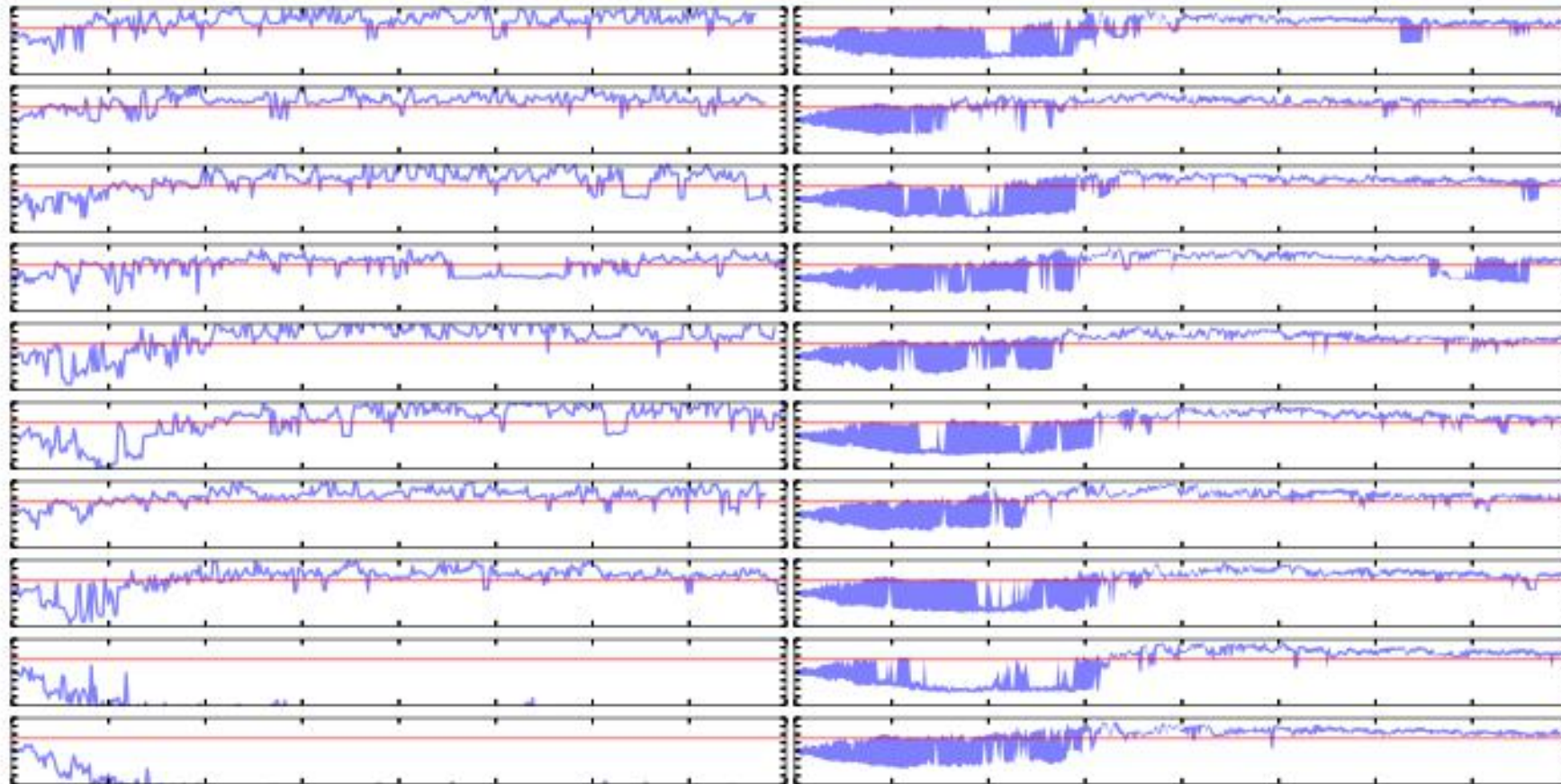




for success. The first experiment uses a simulated swimming snake (see Figure 2), which receives a reward equal to its speed along the  $x$ -axis, either forward or backward. However, once the swimmer swims far enough forward, it crosses a “finish line” and receives a larger reward. Therefore, the best learning strategy is to explore in both directions until the bonus reward is discovered, and then commit to swimming forward. As illustrated in Figure 6 in

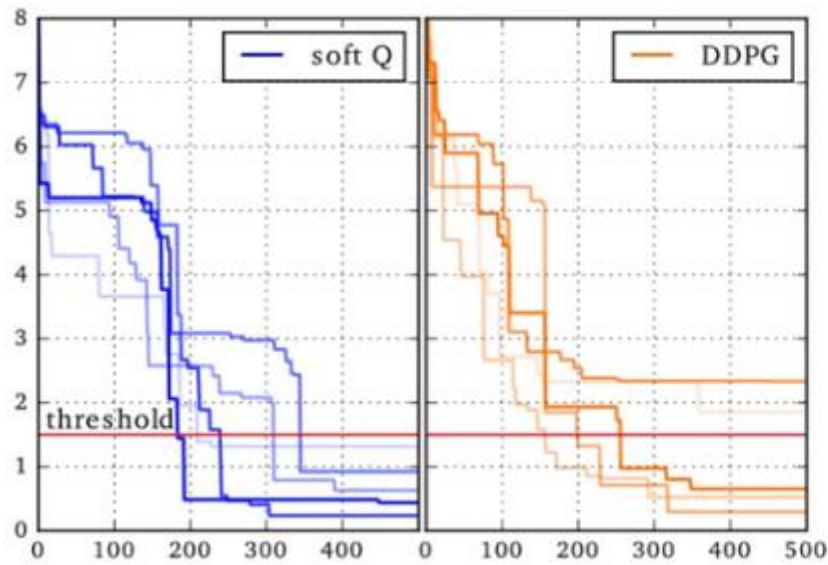


(a) Swimming snake

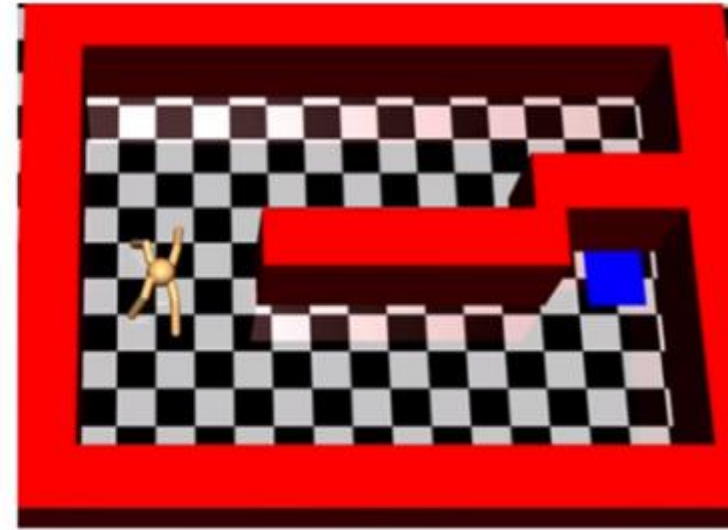


(a) DDPG

(b) Soft Q-learning



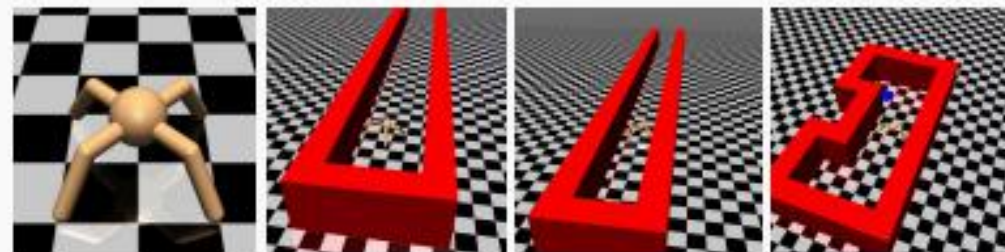
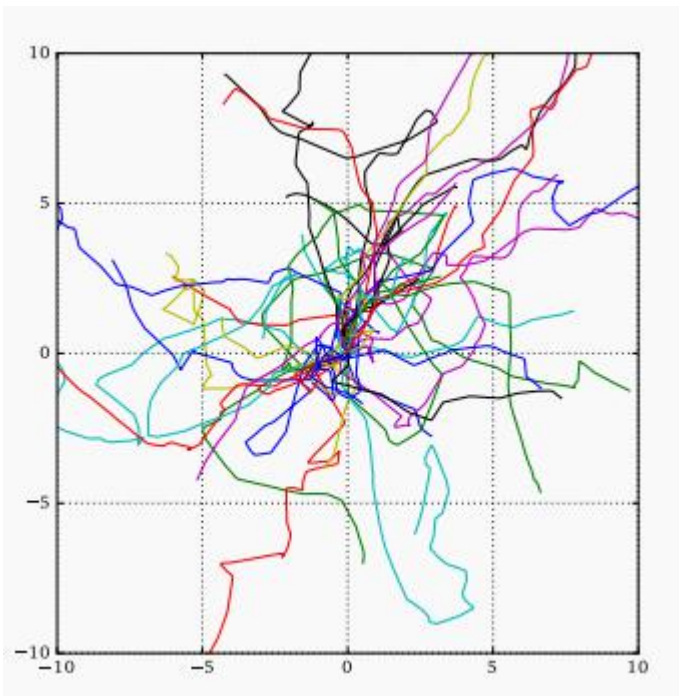
(b) Quadruped (lower is better)



(b) Quadrupedal robot

the performance of DDPG and our method. The curves show the minimum distance to the target achieved so far and the threshold equals the minimum possible distance if the robot chooses the upper passage. Therefore, successful exploration means reaching below the threshold. All policies trained with our method manage to succeed, while only 60% policies trained with DDPG converge to choosing the lower passage.





(a)

(b)

(c)

(d)

Figure 4. Quadrupedal robot (a) was trained to walk in random directions in an empty pretraining environment (details in Figure 7, see Appendix D.3), and then finetuned on a variety of tasks, including a wide (b), narrow (c), and U-shaped hallway (d).

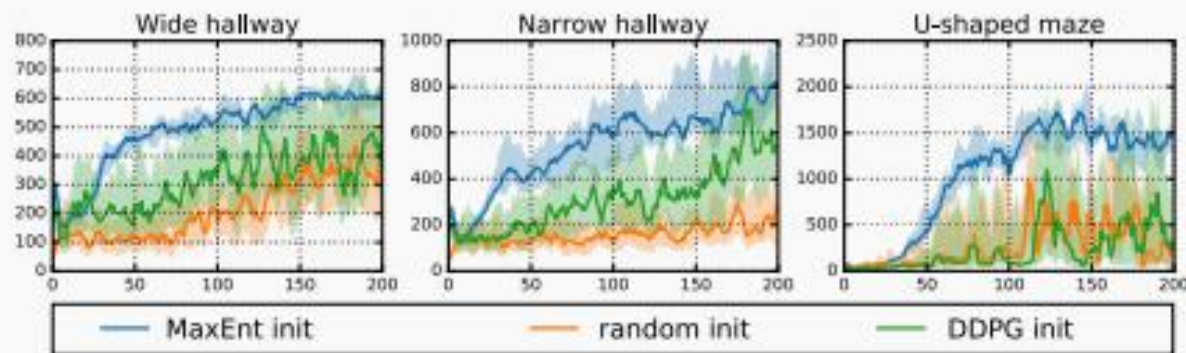


Figure 5. Performance in the downstream task with fine-tuning (MaxEnt) or training from scratch (DDPG). The  $x$ -axis shows the training iterations. The  $y$ -axis shows the average discounted return. Solid lines are average values over 10 random seeds. Shaded regions correspond to one standard deviation.