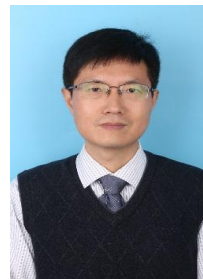


Learning Robust Policy against Disturbance in Transition Dynamics via State-Conservative Policy Optimization



Yufei Kuang, Miao Lu, Jie Wang*, Qi Zhou, Bin Li,
Houqiang Li
University of Science and Technology of China, China

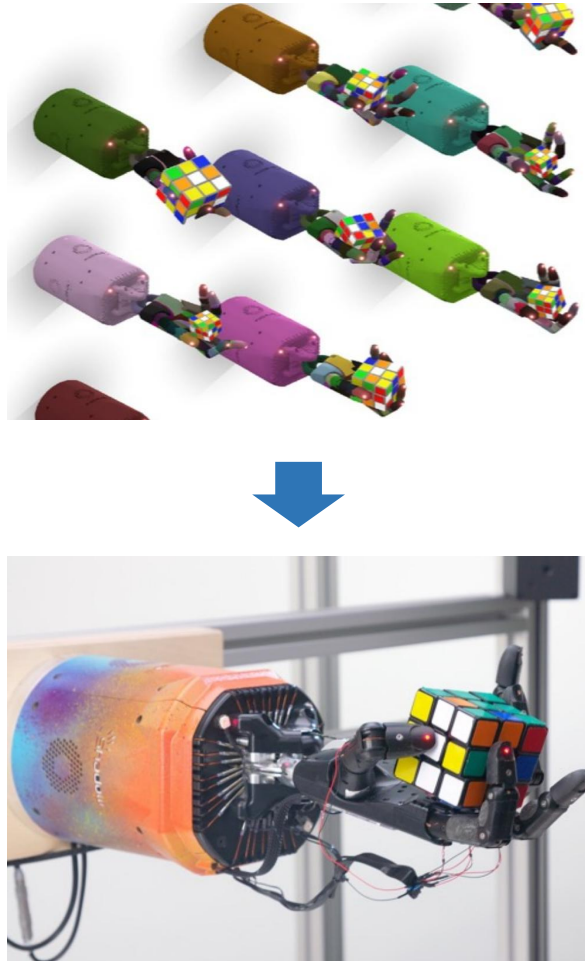
* Corresponding author



- 1** Introduction
- 2** Method
- 3** Experiments
- 4** Conclusion

Part 1

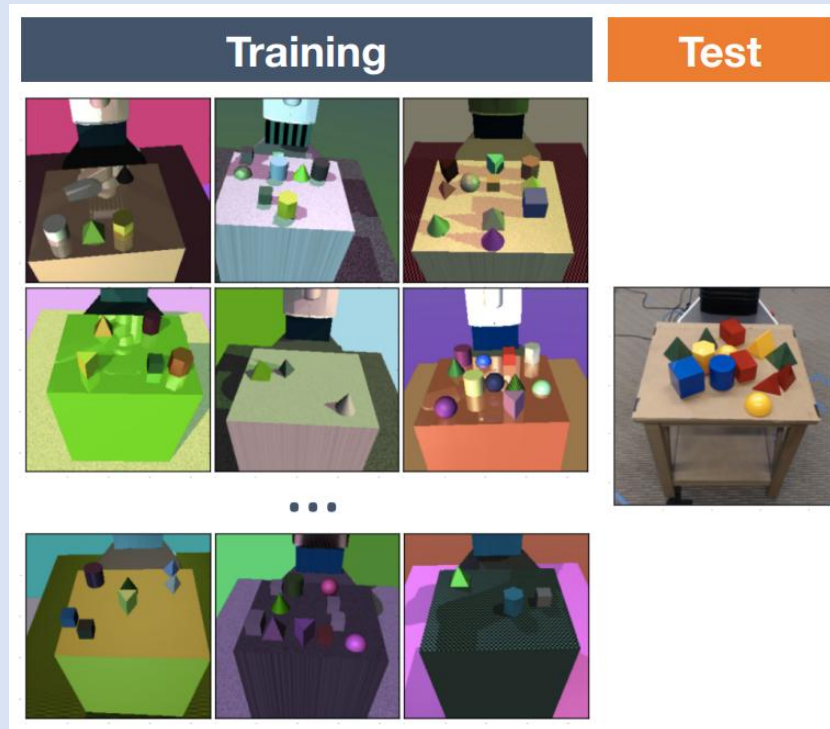
Introduction



- DRL algorithms can perform poorly in real-world tasks due to
 - test-generalization
 - simulation-transfer
- Example:
 - policies trained to control robots in simulators can fail in environments with different mass or friction

domain randomization

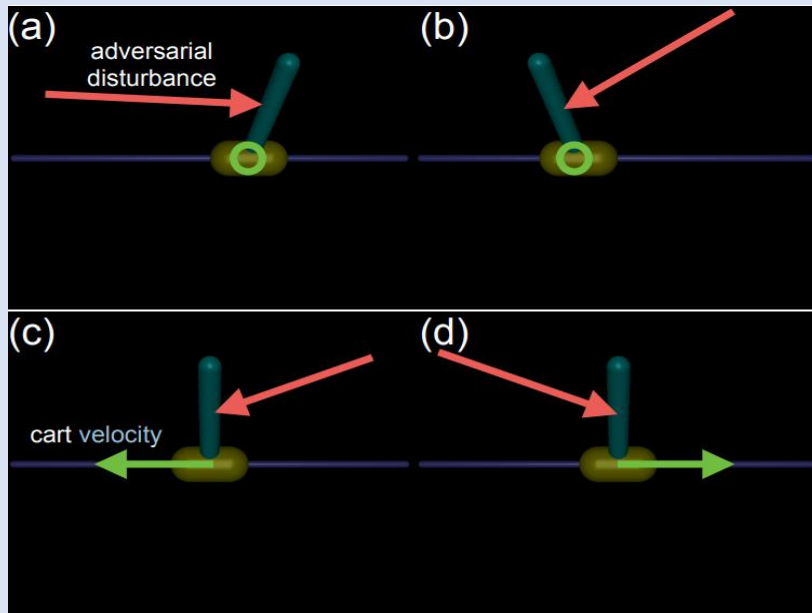
$$J_{\mathcal{P}}(\pi) \triangleq \mathbb{E}_{p \sim \mathcal{P}, \pi} \left[\sum_{t=0}^{+\infty} \gamma^t r(s_t, a_t) \right]$$



- Domain randomization:
 - uses different dynamics to approximate disturbance in target environments.
- Limitation:
 - requires pre-given uncertainty sets

robust adversarial RL

$$J_{\mathcal{P}}(\pi) \triangleq \inf_{p \in \mathcal{P}} \mathbb{E}_{p, \pi} \left[\sum_{t=0}^{+\infty} \gamma^t r(s_t, a_t) \right]$$



- Robust adversarial reinforcement learning:
 - models the disturbance as trainable forces and learns by adversarial training.
- Limitation:
 - requires extra forces manually designed for each task

Challenges:

1. require task-specific prior knowledge to model the disturbance
2. assume control of specially designed simulators

Our goals:

learn robust policies *without modeling the disturbance in advance*

Part 2

Method

Intuition: any disturbance in transition dynamics eventually influences the process via the change of future states

$$\mathcal{T}_{\mathcal{P}}^{\pi} Q \triangleq r(s, a) + \gamma \inf_{p(\cdot|s,a) \in \mathcal{P}(s,a)} \mathbb{E}_{s' \sim p(\cdot|s,a), a' \sim \pi(\cdot|s')} [Q(s', a')] \quad (1)$$

↓ If \mathcal{P} is bounded by Wasserstein distance, then the dual problem of (1) is (2)

$$\mathcal{T}_{\mathcal{P}_{\epsilon}}^{\pi} Q(s, a) = r(s, a) + \gamma \sup_{\lambda \geq 0} \mathbb{E}_{\tilde{s} \sim p_0(\cdot|s,a)} \left\{ \inf_{s' \in \mathcal{S}} \mathbb{E}_{a' \sim \pi(\cdot|s')} [Q(s', a')] + \lambda (d(s', \tilde{s})^p - \epsilon) \right\}. \quad (2)$$

↓ If $V(s)$ is convex and the environment dynamics are deterministic, then (2) is equal to (3)

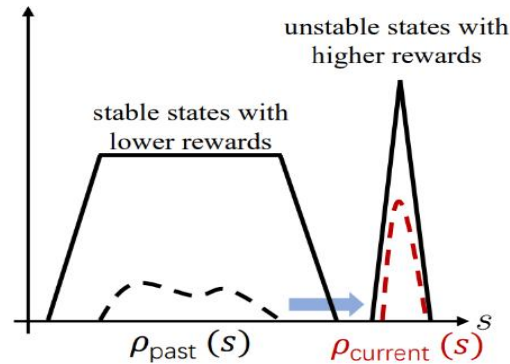
$$\widehat{\mathcal{T}}_{\mathcal{P}_{\epsilon}}^{\pi} Q(s_t, a_t) = r(s_t, a_t) + \gamma \inf_{s' \in B_{\epsilon}(s_{t+1})} \mathbb{E}_{a' \sim \pi(\cdot|s')} [Q(s', a')]. \quad (3)$$

reduces the constrained optimization problem in transition dynamics to a simpler one in state space

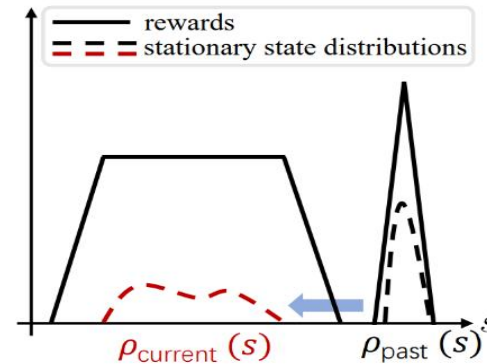
State-Conservative Markov Decision Process (SC-MDP):

Objective:

$$\begin{aligned}
 J_{\epsilon-\mathcal{S}}(\pi) &\triangleq \mathbb{E}_{s_0 \sim d_0} \left[\inf_{s'_0 \in B_\epsilon(s)} \mathbb{E}_{a'_0 \sim \pi(\cdot | s'_0)} [r(s'_0, a'_0)] \right. \\
 &+ \gamma \mathbb{E}_{s_1 \sim p_0(\cdot | s'_0, a'_0)} \left[\inf_{s'_1 \in B_\epsilon(s_1)} \mathbb{E}_{a'_1 \sim \pi(\cdot | s'_1)} [r(s'_1, a'_1)] \right. \\
 &\left. \left. + \gamma \mathbb{E}_{s_2 \sim p_0(\cdot | s'_1, a'_1)} \left[\inf_{s'_2 \in B_\epsilon(s_2)} \mathbb{E}_{a'_2 \sim \pi(\cdot | s'_2)} [r(s'_2, a'_2) + \dots] \dots \right] \right] \right]
 \end{aligned}$$



(a) An example MDP.



(b) Illustrate the SC-MDP.

An illustration of the state-conservative MDP

Based on the above definition of SC-MDP, we can define the corresponding Q function, which is exactly the fixed point of \mathbb{T}_ϵ^π

$$\mathbb{T}_\epsilon^\pi Q(s, a) = r(s, a) + \gamma \mathbb{E}_{\tilde{s} \sim p_0(\cdot|s, a)} \left[\inf_{s' \in B_\epsilon(\tilde{s})} \mathbb{E}_{a' \sim \pi(\cdot|s')} [Q(s', a')] \right]$$

Thus, the traditional strategy iterative algorithm is still convergent in the conservative state Markov decision-making process (Algorithm1) :

Algorithm 1 State-Conservative Policy Iteration

- 1: **Input:** A given MDP with transition p_0 , an initial policy π^0 , a non-negative real ϵ and iteration step K .
 - 2: **for** $k = 0, 1, \dots, K$ **do**
 - 3: Perform **State-Conservative Policy Evaluation** till convergence to obtain $Q_{\epsilon-\mathcal{S}}^{\pi^k}$.
 - 4: Perform **Policy Improvement** step by $\pi^{k+1}(\cdot|s) \leftarrow \arg \max_{\pi \in \Delta_{\mathcal{A}}} \inf_{s' \in B_\epsilon(s)} \mathbb{E}_{a' \sim \pi} [Q_{\epsilon-\mathcal{S}}^{\pi^k}(s', a')]$.
 - 5: **end for**
 - 6: **Output:** an estimation π^{K+1} of the optimal policy.
-

We extend the SC-PI algorithm to a model-free actor-critic algorithm: state-conservative policy optimization (SCPO).

Policy Evaluation (Train Critic)

where:

$$J_{\epsilon-Q}(\theta) \triangleq \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} \left[\frac{1}{2} \left(Q_{\theta}(s_t, a_t) - \hat{Q}(s_t, a_t) \right)^2 \right]$$
$$\hat{Q}(s_t, a_t) \triangleq r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim p_0} [\hat{V}(s_{t+1})]$$
$$\hat{V}(s_{t+1}) \triangleq \inf_{s' \in B_f(s_{t+1})} \mathbb{E}_{a' \sim \pi_{\phi}(\cdot | s')} [Q_{\theta}(s', a')]$$

Policy Improvement (Train Actor)

$$J_{\epsilon-\pi}(\phi) \triangleq \mathbb{E}_{S_t \sim D} \left[\inf_{s \in B_{\epsilon}(s_t)} \mathbb{E}_{a \sim \pi_{\phi}(\cdot | s)} [Q_{\theta}(s, a)] \right].$$

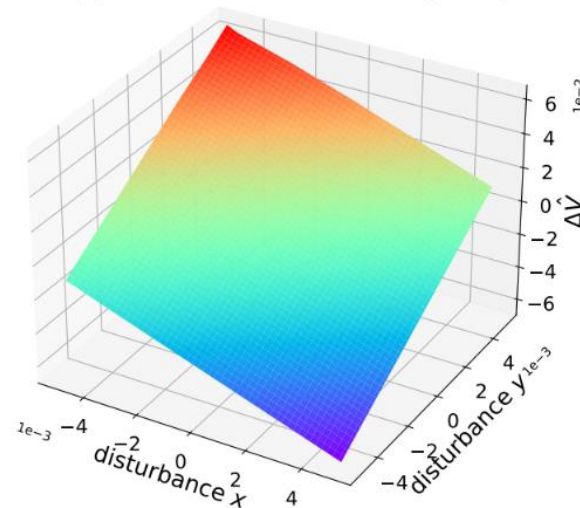


How to solve the following problem

$$\inf_{s' \in B_{\epsilon}(s)} \mathbb{E}_{a' \sim \pi_{\phi}(\cdot | s')} [Q_{\theta}(s', a')]$$

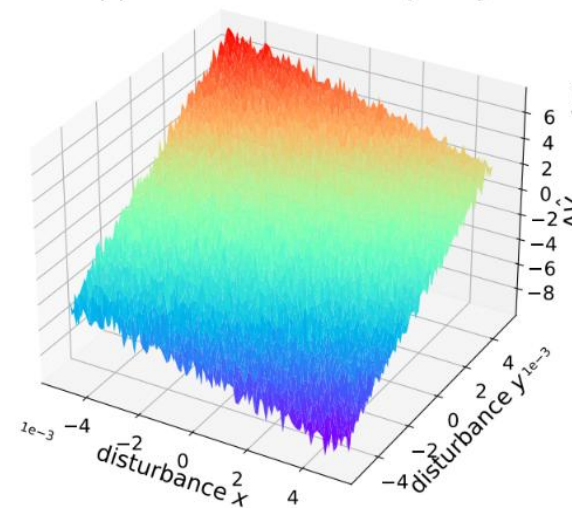
We use a gradient based method that approximately solves the constrained optimization problem efficiently.

Hopper with deterministic policy



(a) Deterministic policy.

Hopper with stochastic policy



(b) Stochastic policy.

$$\inf_{s' \in B_\epsilon(s)} \bar{U}_{\theta, \phi, s}(s') = U_{\theta, \phi}(s) - \epsilon \|\nabla_s U_{\theta, \phi}(s)\|_1$$

where $U_{\theta, \phi}(s) \triangleq \mathbb{E}_{a \sim \pi_\phi(\cdot|s)} [Q_\theta(s, a)]$

Algorithm 2: State-Conservative Soft Actor-Critic

```
1: Input: Critic  $Q_{\theta_1}, Q_{\theta_2}$ . Actor  $\pi_\phi$ . Initial temperature parameter  $\alpha$ . Step size  $\beta_Q, \beta_\pi, \beta_\alpha$ . Target smoothing coefficient  $\tau$ .
2:  $\bar{\theta}_1 \leftarrow \theta_1, \bar{\theta}_2 \leftarrow \theta_2, \mathcal{D} \leftarrow \emptyset$ .
3: for each iteration do
4:   for each environment step do
5:      $a_t \sim \pi(\cdot|s_t), s_{t+1} \sim p_0(\cdot|s_t, a_t)$ .
6:      $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t, r(s_t, a_t), s_{t+1}\}$ .
7:   end for
8:   for each training step do
9:      $\theta_i \leftarrow \theta_i - \beta_Q \nabla J_{\epsilon-Q}(\theta)$  for  $i = 1, 2$ .
10:     $\phi \leftarrow \phi - \beta_\pi \nabla J_{\epsilon-\pi}(\phi)$ .
11:     $\alpha \leftarrow \alpha - \beta_\alpha \nabla J_\alpha(\alpha)$ .
12:     $\bar{\theta}_i \leftarrow \tau\theta_i + (1 - \tau)\bar{\theta}_i$  for  $i = 1, 2$ .
13:   end for
14: end for
15: Output:  $\theta_1, \theta_2, \phi$ 
```

strengths:

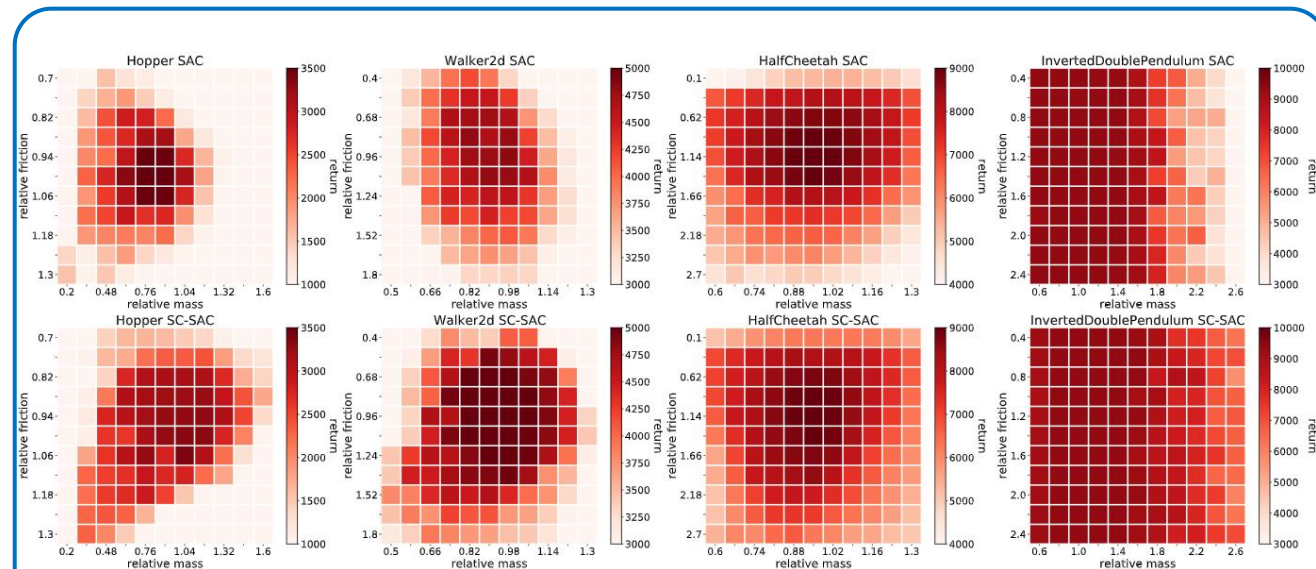
- simple to implement
- does not require task-specific

prior knowledge to model the

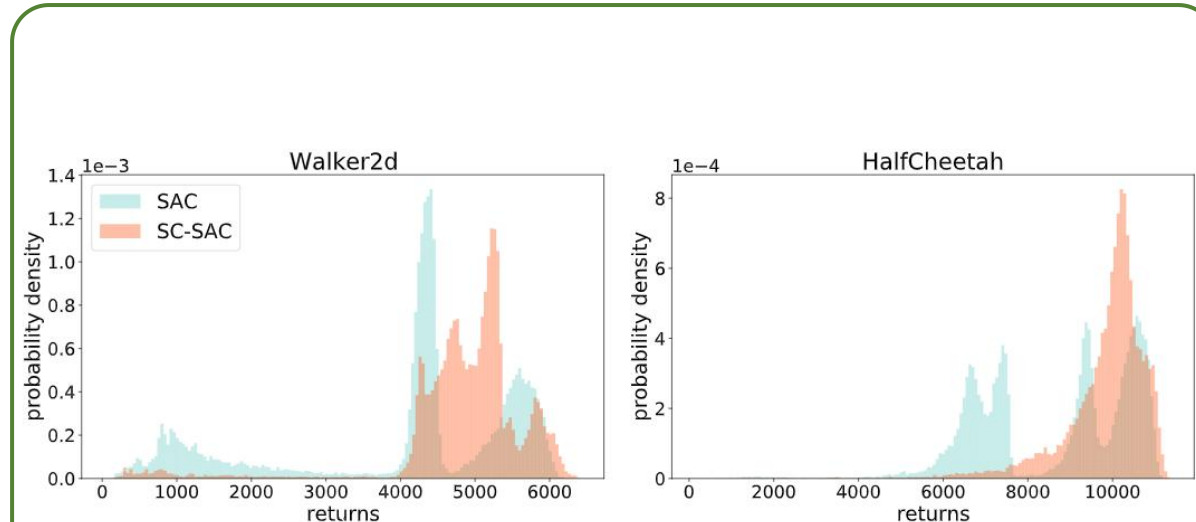
disturbance.

Part 3 Experiments

Results (1): comparison with the original SAC algorithm



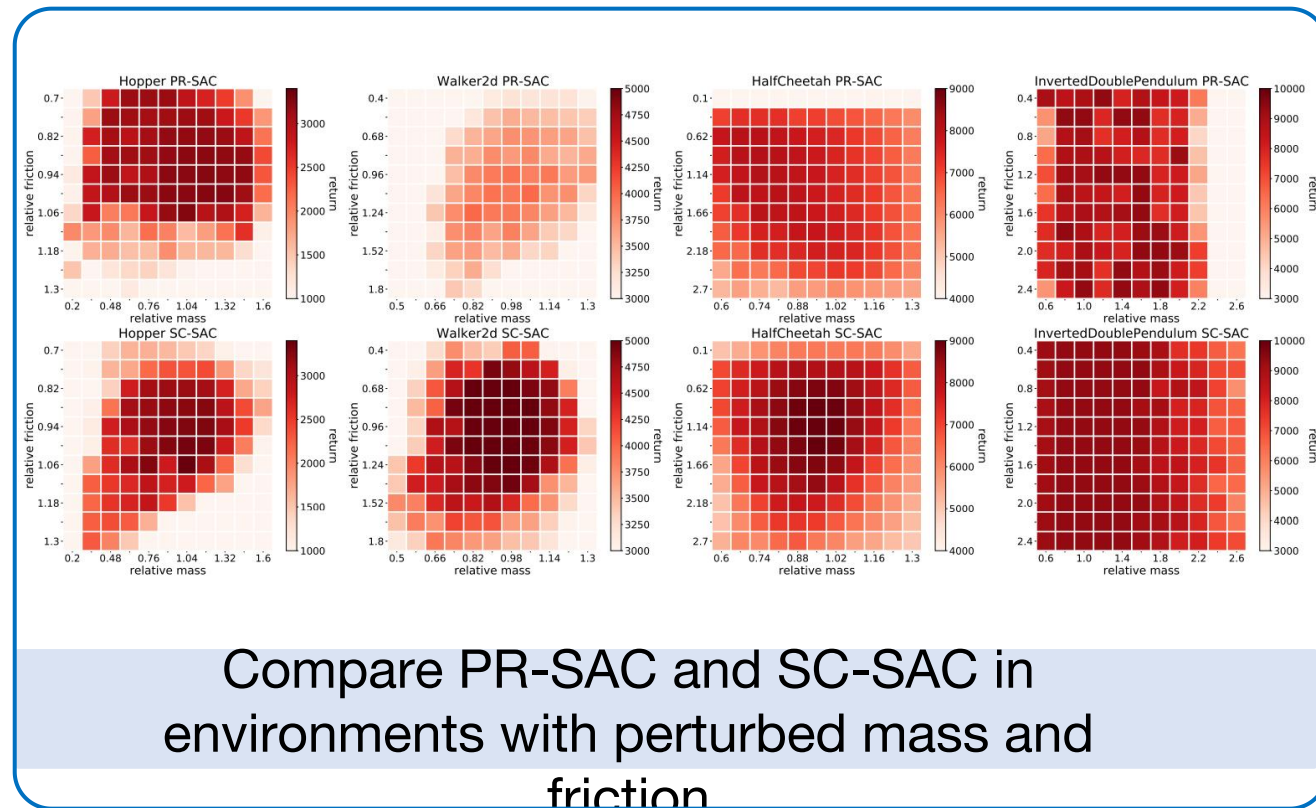
Compare SAC and SC-SAC in environments with perturbed mass and friction



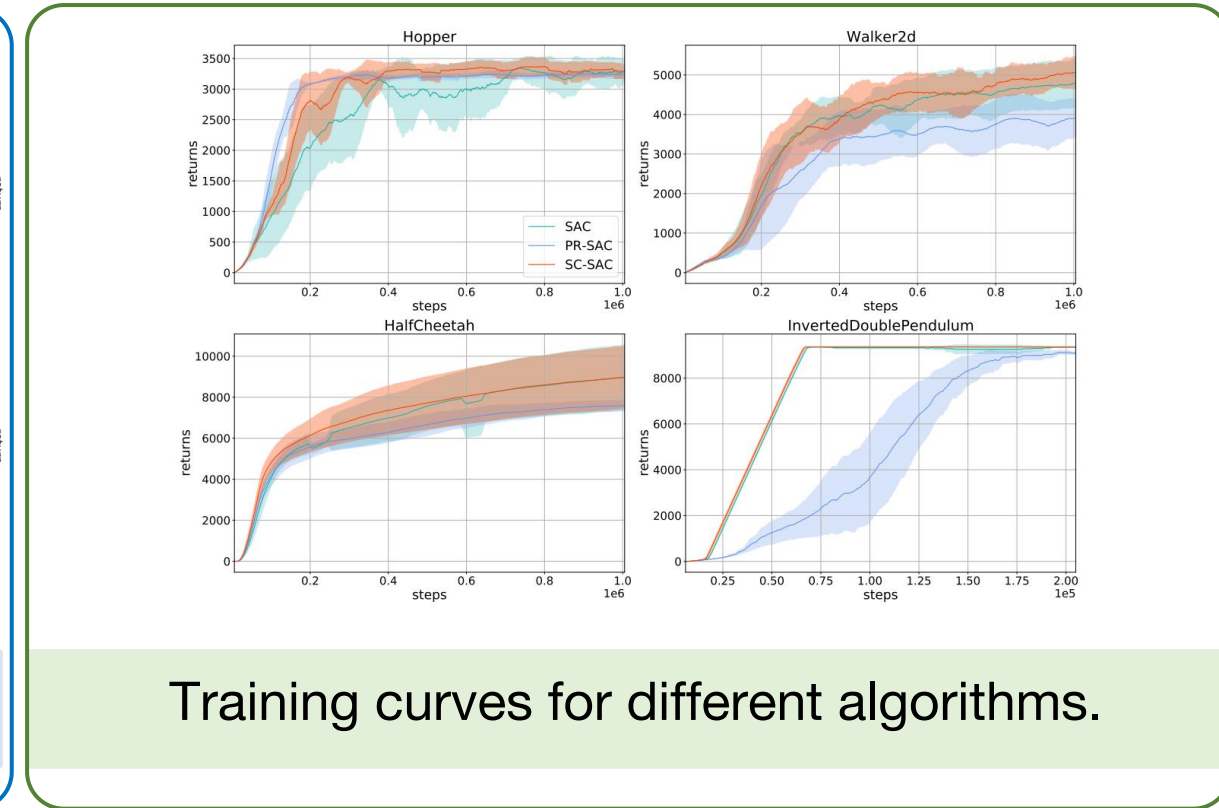
Compare the return distribution between SAC and SC-SAC with truncated Gaussian parameters

Conclusion: SC-SAC is more robust than SAC when generalize to target environments

Results (2): comparison with the action robust algorithm PR-MDP



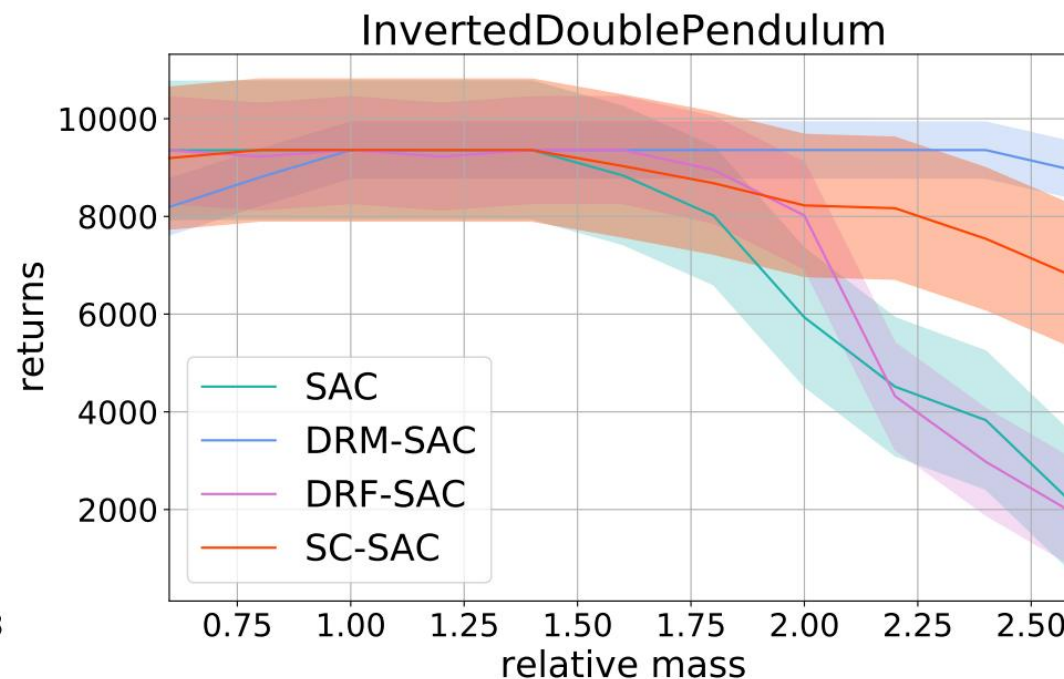
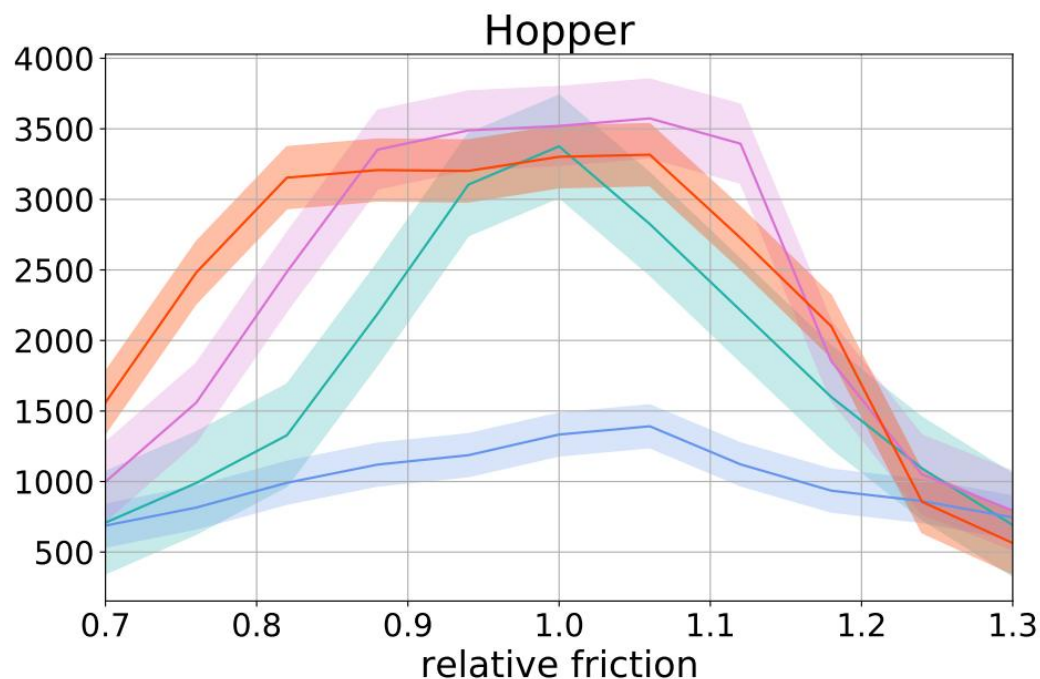
Compare PR-SAC and SC-SAC in environments with perturbed mass and friction



Training curves for different algorithms.

Conclusion: SC-SAC achieves higher average returns than PR-SAC in most tasks

Results (3): comparison with domain randomization (DR)



Compare SCPO with DR in target environments with one parameter perturbed.

Conclusion: though the perturbations on both mass and friction are unmodeled during training, SC-SAC trained policies are robust to them consistently

Part 4

Conclusion

- A novel algorithm SCPO to learn robust policies
 - *simple to implement* and apply to existing actor-critic algorithms
 - does not require *task-specific knowledge* to model the disturbance
 - *learns robust policies* against disturbance in practice

See our Homepage for More Details

<https://miralab.ai/>



<https://miralab.ai/people/yufei-kuang/>



**MANY
THANKS !**