

Q-Pensieve: Boosting Sample Efficiency of Multi-Objective RL Through Memory Sharing of Q-Snapshots

Wei Hung^{12*}, Bo-Kai Huang^{1*}, Ping-Chun Hsieh¹, Xi Liu³

ICLR 2023

1. Department of Computer Science, National Yang Ming Chiao Tung University

2. Research Center for Information Technology Innovation, Academia Sinica

3. Applied Machine Learning, Meta AI

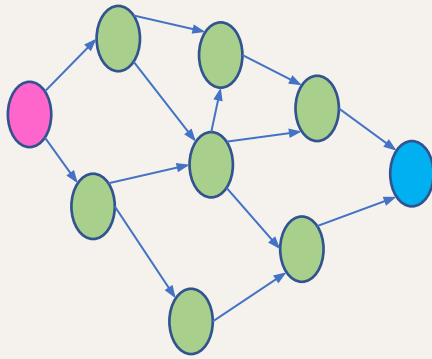


Introduction

- We identify the critical sample inefficiency issue in the existing MORL algorithms for continuous control
- We propose Q-Pensieve, which is a policy improvement scheme for enhancing the data sharing capability across policies
- We substantiate the concept of Q-Pensieve policy iteration by proposing the technique of Q replay buffer and arrive at a practical actor-critic type practical implementation

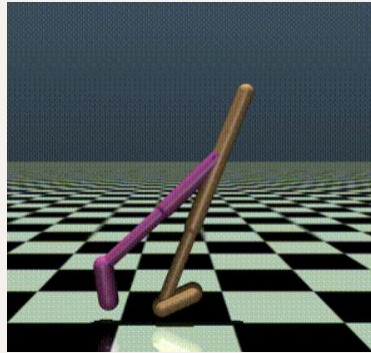
Multi-Objective Applications in Real World

Communication networks



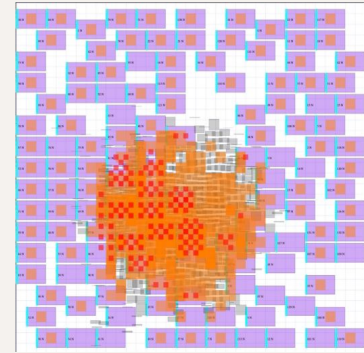
- Throughput and latency

Robotics



- Speed, survival bonus and control cost

Chip Design



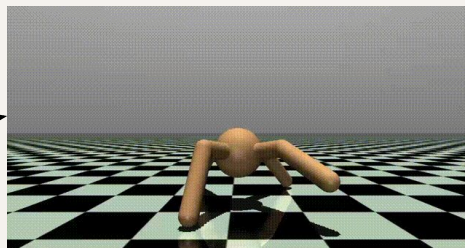
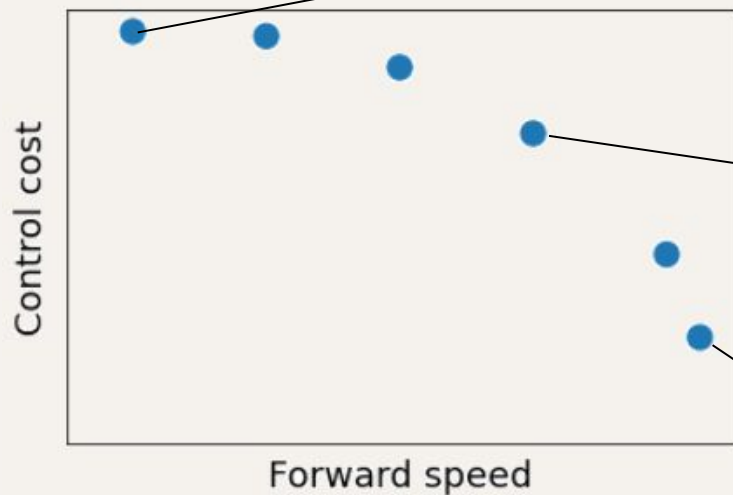
- Wirelength, routing congestion, and density

(Credit by Mirhoseini et al.)

MORL Learns Diverse Behaviors

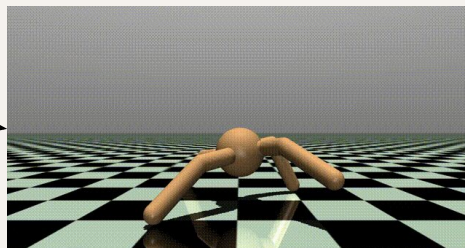
Forward speed: running speed

Control cost: -energy consumption



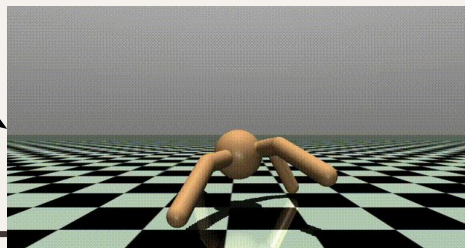
$\lambda=[0.1,0.9]$

Balance with minimal control cost



$\lambda=[0.5,0.5]$

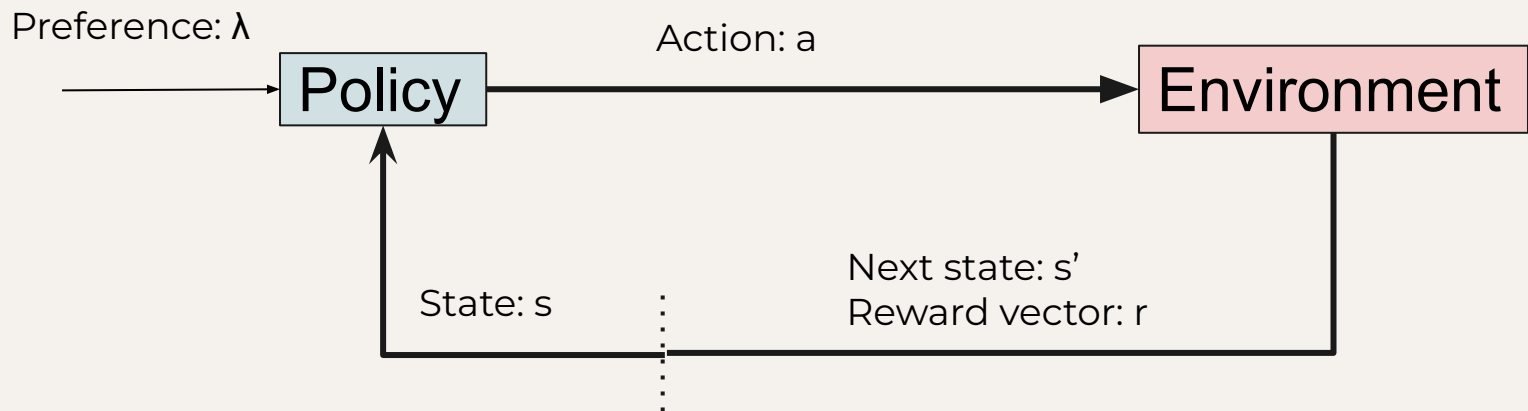
Walk gracefully



$\lambda=[0.9,0.1]$

Sprint at all cost!!

MORL Formulation



$$Q_{\pi_{\theta}}(s, a) := \mathbf{E}[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}_t | s_0 = s, a_0 = a; \pi_{\theta}] \quad (\text{Q-functions in MORL})$$

$$J(\pi_{\theta}) := \mathbf{E}_{s_0, a_0 \sim \pi_{\theta}} [Q_{\pi_{\theta}}(s_0, a_0)] \quad (\text{Vector-Valued Total Return})$$

Goal: Given all λ , learn $\boldsymbol{\pi}_{\theta}$ that maximizes utility $\lambda^{\top} J(\pi_{\theta})$

Convex Coverage Set (CCS)

● : Non-CCS

It is a total return vector of some π

$$J(\pi) = \sum_t \gamma^t \mathbf{r}(s, a)$$

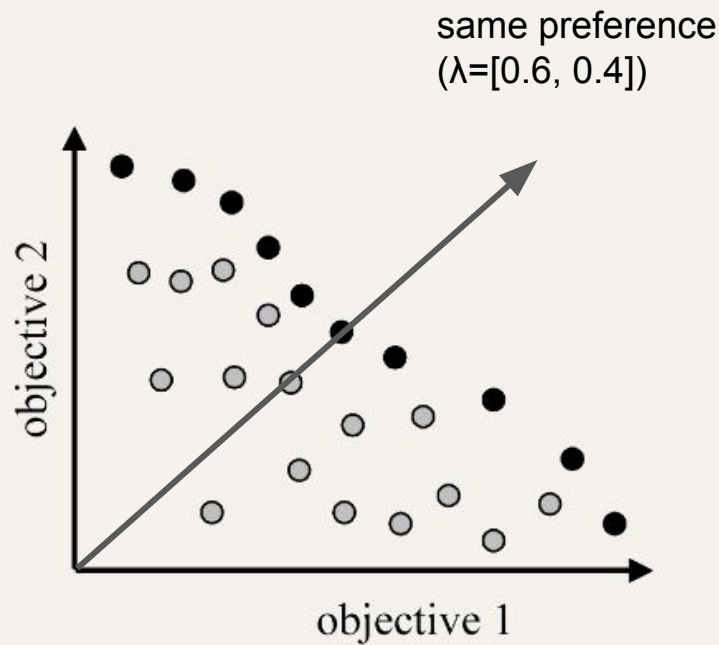
● : CCS

It is an optimal return vector for some λ

$$\lambda^\top J(\pi^*) \geq \lambda^\top J(\pi), \forall \pi$$

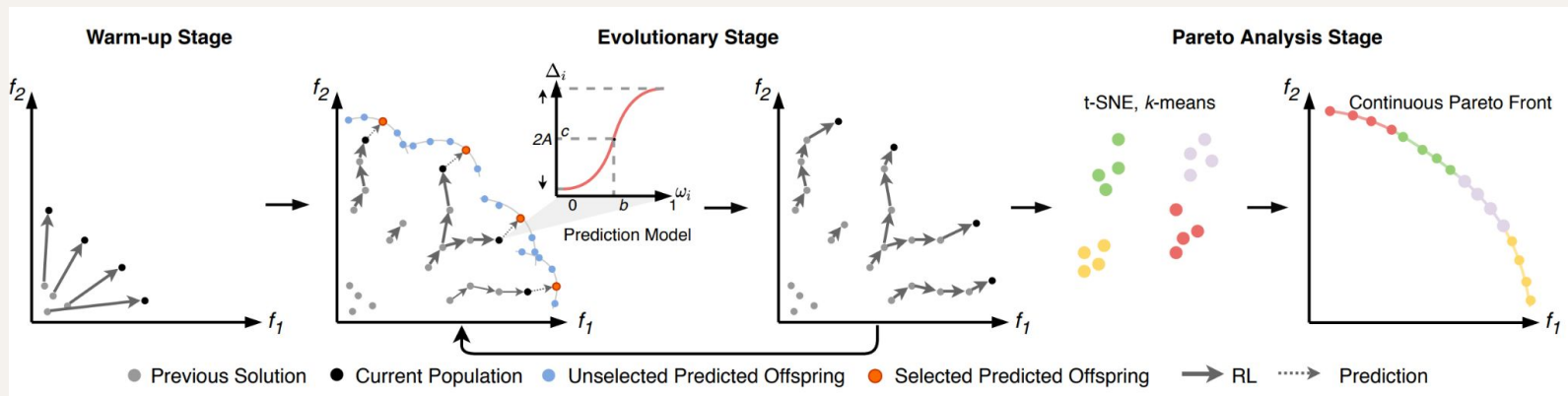
Different Aspects:

- Maximize hypervolume (Area covered by reward vector)
- Maximize $\lambda^\top J(\pi), \forall \lambda$



Existing Solutions to MORL

- Explicit search - PGMORL
 - Evolutionary search for CCS
 - Issue: Sample inefficiency!

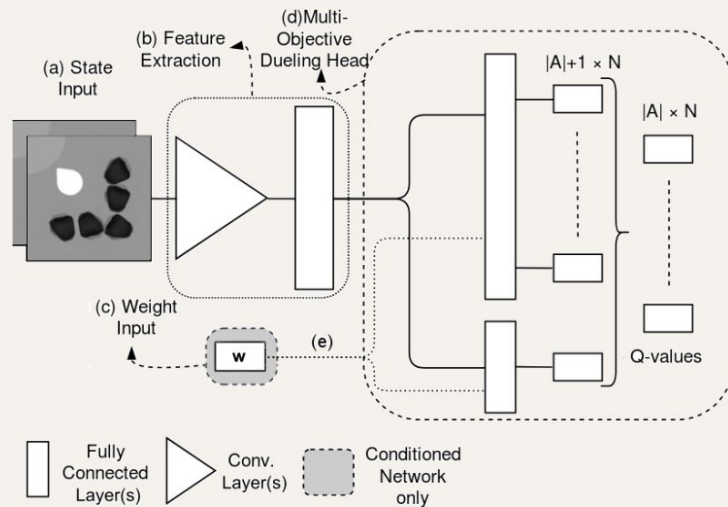


[Xu et al.,2020]

1. Jie Xu, Yunsheng Tian, Pingchuan Ma, Daniela Rus, Shinjiro Sueda, and Wojciech Matusik. Prediction-Guided Multi-Objective Reinforcement Learning for Continuous Robot Control, ICML 2020

Existing Solutions to MORL

- Implicit search - Conditioned Networks (CN)
 - Preference-dependent MO-DQN
 - Issue: No policy improvement guarantee!



[Abels et al., 2019]

1. Axel Abels, Diederik M. Roijers, Tom Lenaerts, Ann Nowé, and Denis Steckelmacher. Dynamic Weights in Multi-Objective Deep Reinforcement Learning, ICML 2019

Envelope Q-Learning

- Align one preference with optimal rewards that may have been explored under other preferences
- Bellman backup operator

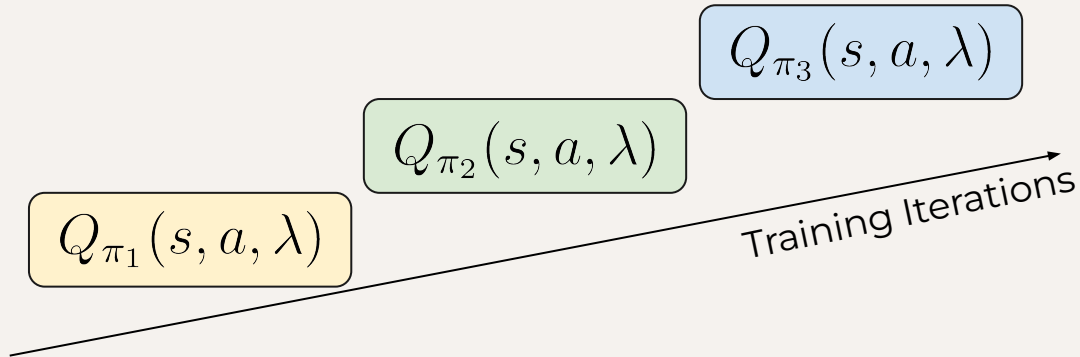
$$\mathcal{T}_\pi Q(s, a; \lambda) = \mathbf{r}(s, a) + \gamma \mathbf{E}_{(s', a') \sim (\mathcal{P}, \pi)} Q(s', a'; \lambda)$$

- Optimality filter for multi-objective Q

$$(\mathcal{H}Q)(s; \lambda) = \arg_Q \sup_{\underline{a \in \mathcal{A}, \lambda' \in \Lambda}} \lambda^\top Q(s, a; \lambda')$$

1. Runzhe Yang, Xingyuan Sun, and Karthik Narasimhan. A generalized algorithm for multi-objective reinforcement learning and policy adaptation, NeurIPS 2019

Our Idea: Memory Sharing of Snapshots (Called Q-Pensieve)



- Policy-level knowledge sharing: Snapshots can boost the learning in future iterations
- Each Q-network $Q_{\pi_k}(s, a, \lambda)$ can be good for some preference vector λ



Pensieve: A magical device used to review and store memories

Q-Pensieve Policy Improvement Update

MO Soft Policy improvement

$$\pi_{k+1}(\cdot | \cdot; \lambda) = \arg \min_{\pi' \in \tilde{\Pi}} D_{\text{KL}} \left(\pi'(\cdot | s; \lambda) \parallel \frac{\exp \left(\sup_{\lambda' \in W_k(\lambda), \mathbf{Q}' \in \mathcal{Q}_k} \left(\frac{1}{\alpha} \lambda^\top \mathbf{Q}'(\cdot, \cdot; \lambda') \right) \right)}{Z_{\mathcal{Q}_k}(s)} \right)$$

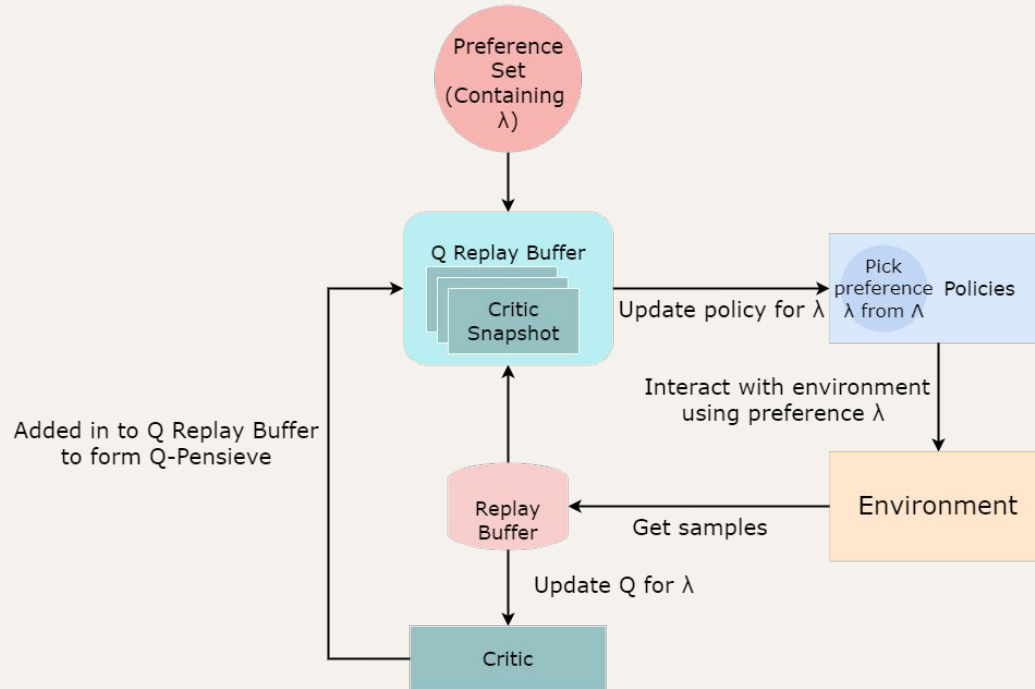
$W_k(\lambda)$ is a subset of Λ containing λ

$\mathcal{Q}_k(\lambda)$ is a set of Q-snapshots containing $Q_{\pi_k}(s, a, \lambda)$

Theorem (Convergence of Q-Pensieve)

Repeated application of soft policy evaluation and soft improvement to any $\pi \in \Pi$ converges to a policy π^* such that $\lambda^\top Q^{\pi^*}(s, a; \lambda) \geq \lambda^\top Q^\pi(s, a; \lambda)$ for all π , and all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $\lambda \in \Lambda$

Implementation of Q-Pensieve

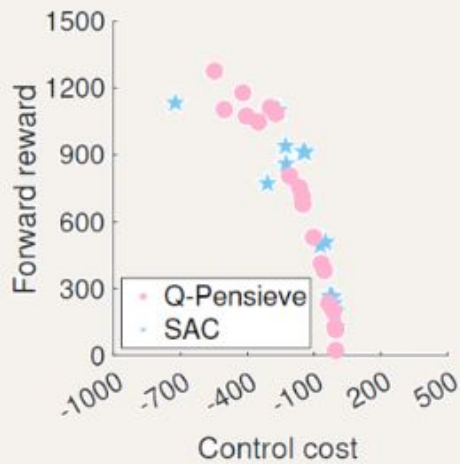


Experimental Results - Comparison with Baselines

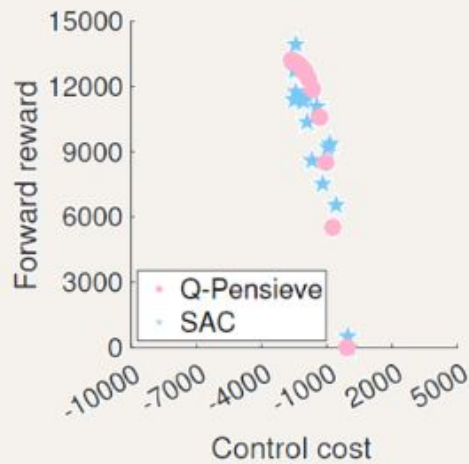
Environments	Metrics	PFA	PFA	PGMORL	PGMORL	CN-DER	Q-Pensieve
		(1.5M steps)	($1.5 \times \beta M$ steps)	(1.5M steps)	($1.5 \times \beta M$ steps)	(1.5M steps)	(1.5M steps)
DST2d	HV($\times 10^2$)	7.43 \pm 3.68	8.67 \pm 1.49	8.10 \pm 1.57	8.13 \pm 1.61	5.36 \pm 4.71	10.21\pm1.40
	UT	-9.27 \pm 6.03	-6.86 \pm 6.06	4.90 \pm 0.44	5.02 \pm 0.35	-5.10 \pm 15.73	7.31\pm0.91
	ED	0.13 \pm 0.11	0.10 \pm 0.08	0.25 \pm 0.18	0.28 \pm 0.18	0.21 \pm 0.17	0.54\pm0.11
HalfCheetah2d	HV($\times 10^7$)	0.73 \pm 0.19	1.31 \pm 0.26	0.53 \pm 0.17	0.28 \pm 0.29	2.08 \pm 0.54	3.82\pm0.27
	UT($\times 10^3$)	0.31 \pm 0.20	1.02 \pm 0.40	-0.28 \pm 0.94	0.09 \pm 0.17	5.09 \pm 3.57	5.61\pm0.31
	ED	0.08 \pm 0.10	0.10 \pm 0.06	0.01 \pm 0.00	0.11 \pm 0.05	0.02 \pm 0.01	0.54\pm0.08
Hopper2d	HV($\times 10^6$)	0.49 \pm 0.46	1.01 \pm 0.62	0.63 \pm 0.48	1.31 \pm 0.48	0.56 \pm 0.16	1.33\pm0.20
	UT($\times 10^2$)	2.89 \pm 1.93	3.50 \pm 1.85	1.94 \pm 2.46	3.70 \pm 1.78	1.42 \pm 1.00	4.08\pm1.10
	ED	0.31 \pm 0.17	0.41 \pm 0.10	0.31 \pm 0.25	0.31 \pm 0.11	0.04 \pm 0.03	0.43\pm0.09
Ant2d	HV($\times 10^6$)	0.17 \pm 0.05	0.77 \pm 0.53	0.14 \pm 0.03	0.13 \pm 0.04	5.03 \pm 3.60	10.01\pm1.86
	UT($\times 10^2$)	-0.06 \pm 0.01	0.14 \pm 0.14	-0.21 \pm 0.15	-0.18 \pm 0.38	3.68 \pm 2.34	14.04\pm3.03
	ED	0.22 \pm 0.03	0.22 \pm 0.02	0.21 \pm 0.02	0.21 \pm 0.03	0.21 \pm 0.08	0.60\pm0.07
Walker2d	HV($\times 10^6$)	0.52 \pm 0.20	1.05 \pm 0.44	0.83 \pm 0.42	1.28\pm0.66	0.42 \pm 0.09	1.12 \pm 0.36
	UT($\times 10^2$)	0.23 \pm 0.13	0.95 \pm 0.55	0.38 \pm 0.24	1.20 \pm 0.67	3.17 \pm 0.53	6.37\pm1.42
	ED	0.32 \pm 0.06	0.37 \pm 0.09	0.30 \pm 0.10	0.34 \pm 0.12	0.21 \pm 0.11	0.48\pm0.10

Experimental Results - Sample Efficiency of Q-Pensieve

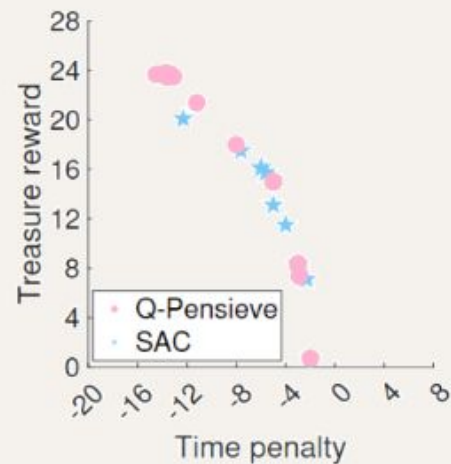
- Select 19 fixed preferences and learn 19 separate SAC models
- We achieve the same performance with only 1/19 of samples used by SAC



(a) Hopper2d



(b) HalfCheetah2d



(c) DST2d

Summary

- We propose Q-Pensieve to boost the sample efficiency of MORL problems
- We present Q-Pensieve soft policy iteration in the tabular setting and show that it preserves the global convergence property
- Our theoretical and experimental results demonstrate that the proposed learning algorithm is indeed a promising approach for MORL problems