

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

Wei Hung<sup>1,2\*</sup>, Bo-Kai Huang<sup>1\*</sup>, Ping-Chun Hsieh<sup>1</sup>, and Xi Liu<sup>3</sup>

<sup>1</sup> Department of Computer Science, National Yang Ming Chiao Tung University, Hsinchu, Taiwan

<sup>2</sup> Research Center for Information Technology Innovation, Academia Sinica, Taipei, Taiwan

<sup>3</sup> Applied Machine Learning, Meta AI, Menlo Park, CA, USA

\* Equal Contribution

Meta



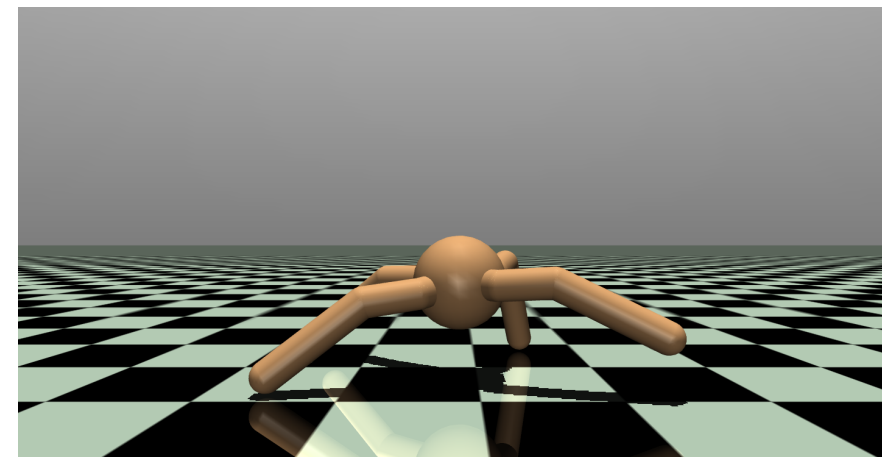
## Introduction

- 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 implementation
- We demonstrate that the proposed **Q-Pensieve** can indeed achieve significantly better empirical sample

## Multi-Objective Reinforcement Learning (MORL)

- Many real-world continuous control problems are also multi-objective tasks by nature, e.g., Control the robot to run fast while consuming as little energy as possible
- The multi-objective RL problems have been extensively studied from two major perspectives
  - ▷ **Explicit search** methods update a policy or a set of policies by explicitly searching for the Pareto front of the reward space (Xu et al., 2020; Kyriakis et al., 2022)
    - It is typically difficult to maintain a sufficiently diverse set of optimal policies for different preferences within a reasonable number of training samples
  - ▷ **Implicit search** methods improve policies for multiple preferences through implicit search (Abels et al., 2019; Yang et al., 2019)
    - Data sharing among policies of different preferences is not guaranteed to achieve policy improvement for all preferences

⇒ We want to propose an efficient way to solve multi-objective tasks.



## Multi-Objective Soft Policy Iteration (MOSPI)

- With Soft Policy Iteration (SPI), we get a robust and good exploration method that can handle continuous actions properly
- Given a preference  $\lambda$ , we generalize the bellman operator to the multi-objective settings
- Multi-Objective Soft Policy evaluation:

$$(\mathcal{T}_{\text{MO}}^{\pi} \mathbf{Q})(\mathbf{s}, \mathbf{a}; \lambda) = r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim \mathcal{P}(\cdot | \mathbf{s}, \mathbf{a}), \mathbf{a}' \sim \pi(\cdot | \mathbf{s}'; \lambda)} [\mathbf{Q}(\mathbf{s}', \mathbf{a}'; \lambda) - \alpha \log \pi(\mathbf{a}' | \mathbf{s}'; \lambda) 1_d]$$

- Multi-Objective Soft Policy improvement:

$$\pi_{k+1}(\cdot, \cdot; \lambda) = \arg \min_{\pi' \in \Pi} D_{\text{KL}}(\pi'(\cdot | \mathbf{s}) || \frac{\exp(\frac{1}{\alpha} \lambda^{\top} \mathbf{Q}^{\pi_k}(\mathbf{s}, \cdot; \lambda))}{Z^{\pi_k}(\mathbf{s})})$$

## Envelope Q-Learning

- Align one preference with optimal rewards that may have been explored under other preferences (Yang et al., 2019)
- Optimality filter for multi-objective  $\mathbf{Q}$ :

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

## Q-Pensieve

- $\mathcal{Q}$  denotes all possible  $\mathbf{Q}$ , then we can achieve policy improvement according to

$$\max_{\mathbf{Q}(\mathbf{s}) \in \mathcal{Q}} \lambda^{\top} \mathbf{Q}$$

- Intuition: Enforce knowledge sharing at the **policy-level** and enhance the sample used in learning a variety of Pareto optimal policies
- Implementation of **Q-Pensieve**: Since we can directly use the  $\mathbf{Q}$  functions of policies in the past without additional training, we store the learned  $\mathbf{Q}$  network in the past iterations as candidates for forming the committee
- **Q Replay Buffer**: The policy update of **Q-Pensieve** would involve both the current  $\mathbf{Q}$ -function and the  $\mathbf{Q}$ -snapshots from the past iterations. We introduce  $\mathbf{Q}$  replay buffer, which could store multiple  $\mathbf{Q}$ -networks in a predetermined manner

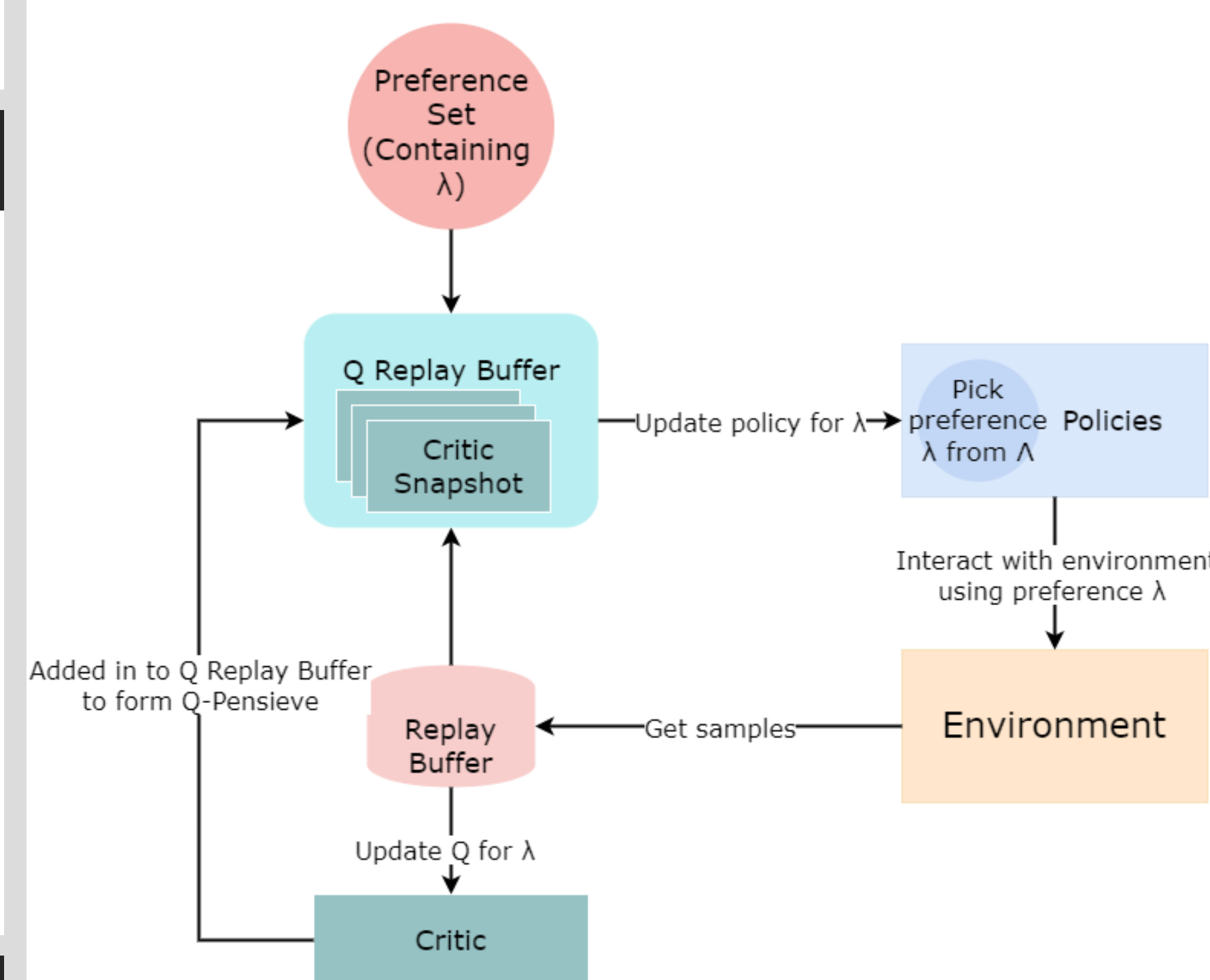
## Q-Pensieve Soft Policy Iteration

- **Q-Pensieve Policy Improvement**: In the policy improvement step of the k-th iteration, for each specific  $\lambda$ , we update the policy as

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

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

## Algorithm



**Q-Pensieve algorithm**: In each iteration, do the following:

- (1) Pick a preference  $\lambda$ , interact with the environment, and add samples to the replay buffer
- (2) Update the critic network by minimizing the critic loss  $\mathcal{L}_{\mathbf{Q}}$
- (3) Build up a preference set  $\mathcal{W}_k(\lambda)$  with  $\lambda$ , sample N preferences and add them to  $\mathcal{W}_k(\lambda)$
- (4) Update the  $\mathbf{Q}$  replay buffer  $\mathcal{Q}_k$  with current critic network
- (5) Update the policy network by minimizing the policy loss  $\mathcal{L}_{\pi}$

Figure: The architecture of Q-Pensieve

$$\mathcal{L}_{\mathbf{Q}}(\phi; \lambda) = \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mu} \left[ \lambda^{\top} (\mathbf{Q}_{\phi}(\mathbf{s}, \mathbf{a}; \lambda) - (r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim \mathcal{P}(\cdot | \mathbf{s}, \mathbf{a})} [V_{\phi}(\mathbf{s}')] ))^2 \right]$$

$$\mathcal{L}_{\pi}(\theta; \lambda) = \mathbb{E}_{\mathbf{s} \sim \mu} \left[ \mathbb{E}_{\mathbf{a} \sim \pi_{\theta}} \left[ \sup_{\lambda' \in \mathcal{W}_k(\lambda), \mathbf{Q} \in \mathcal{Q}_k} \{ \alpha \log (\pi_{\theta}(\mathbf{a} | \mathbf{s}; \lambda)) - \lambda^{\top} \mathbf{Q}_{\phi}(\mathbf{s}, \mathbf{a}; \lambda') \} \right] \right]$$

## Experimental Results - Comparison with Benchmarks

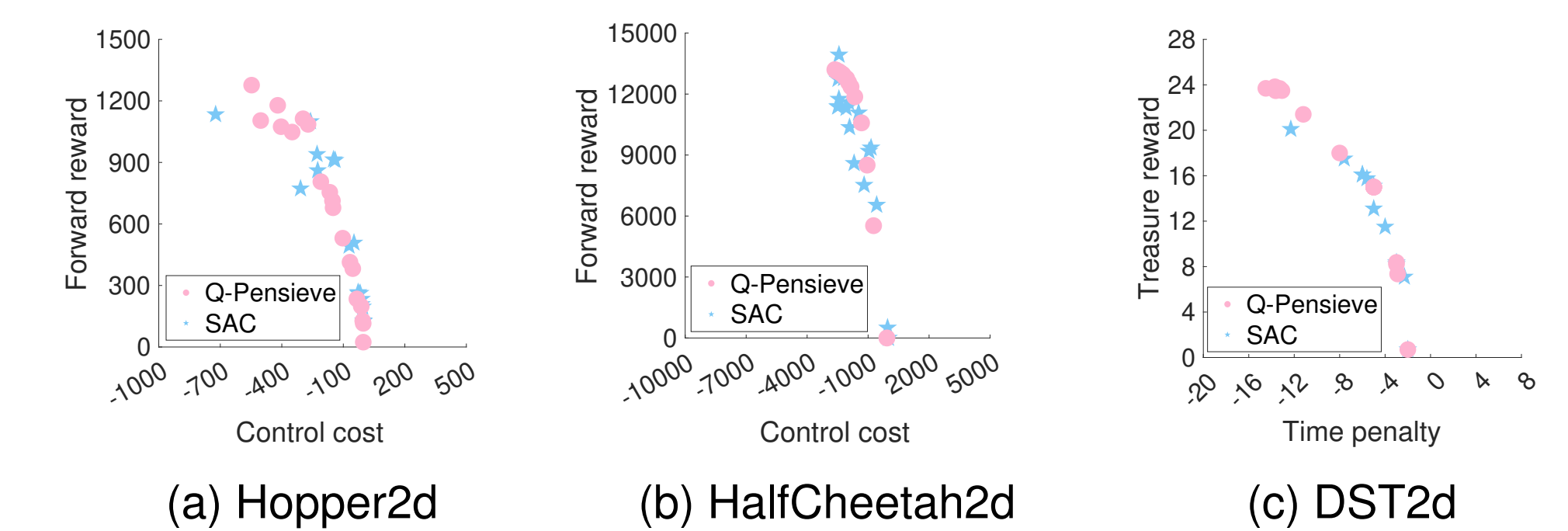
### Performance Metrics

- ▷ HyperVolume (HV): Compute the region dominated by the solution set w.r.t. a given reference point  $\mathbf{r}_0$
- ▷ Utility (UT): Evaluate the performance under linear scalarization

Environments	Metrics	PFA	PGMORL	CN-DER	Q-Pensieve
(Parisi et al., 2014) (Xu et al., 2020) (Abels et al., 2019)					
DST2d	HV( $\times 10^2$ )	7.43	8.10	5.36	<b>10.21</b>
	UT( $\times 10^0$ )	-9.27	4.90	-5.10	<b>7.31</b>
HalfCheetah2d	HV( $\times 10^7$ )	0.73	0.53	2.08	<b>3.82</b>
	UT( $\times 10^3$ )	0.31	-0.28	5.09	<b>5.61</b>
Ant3d	HV( $\times 10^8$ )	-	0.41	13.00	<b>21.87</b>
	UT( $\times 10^3$ )	-	0.18	0.49	<b>1.14</b>
Hopper5d	HV( $\times 10^{13}$ )	-	0.63	3.42	<b>7.24</b>
	UT( $\times 10^2$ )	-	1.48	1.76	<b>3.37</b>

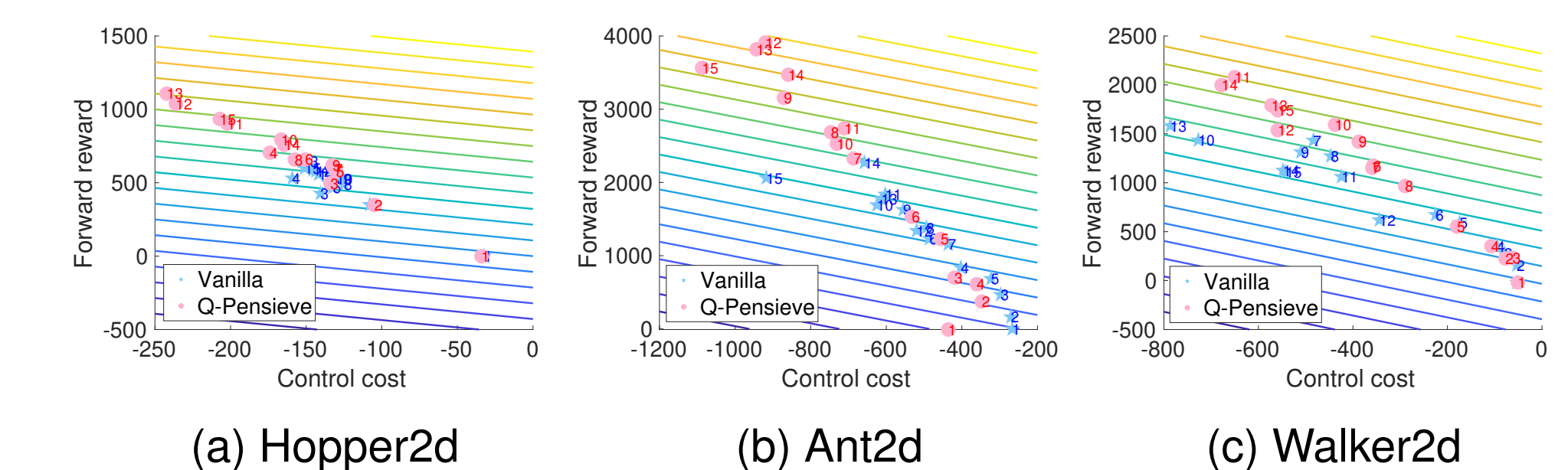
## Experimental Results - Sample Efficiency of Q-Pensieve

We train 19 single-objective SAC models, each for a unique preference, and train **Q-Pensieve** for all preferences. The following figures show the return vectors attained by **Q-Pensieve** and the collection of single-objective SAC models under 19 preferences.



## An Ablation Study on Q Replay Buffer

To verify the effectiveness of the technique of  $\mathbf{Q}$  replay buffer, we compare the performance of **Q-Pensieve** with buffer size equal to 4 and that without using  $\mathbf{Q}$  replay buffer (termed “Vanilla”). The following figures show the return vectors attained under preference  $\lambda = [0.5, 0.5]$  at  $100 \cdot x$  thousand training steps.



## Conclusion

- 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