

Yutong Wang\*, Ke Xue\*, Chao Qian  
(\*Equal contribution)

Email: {wangyt, xuek, qianc}@lamda.nju.edu.cn

## Background and Motivation

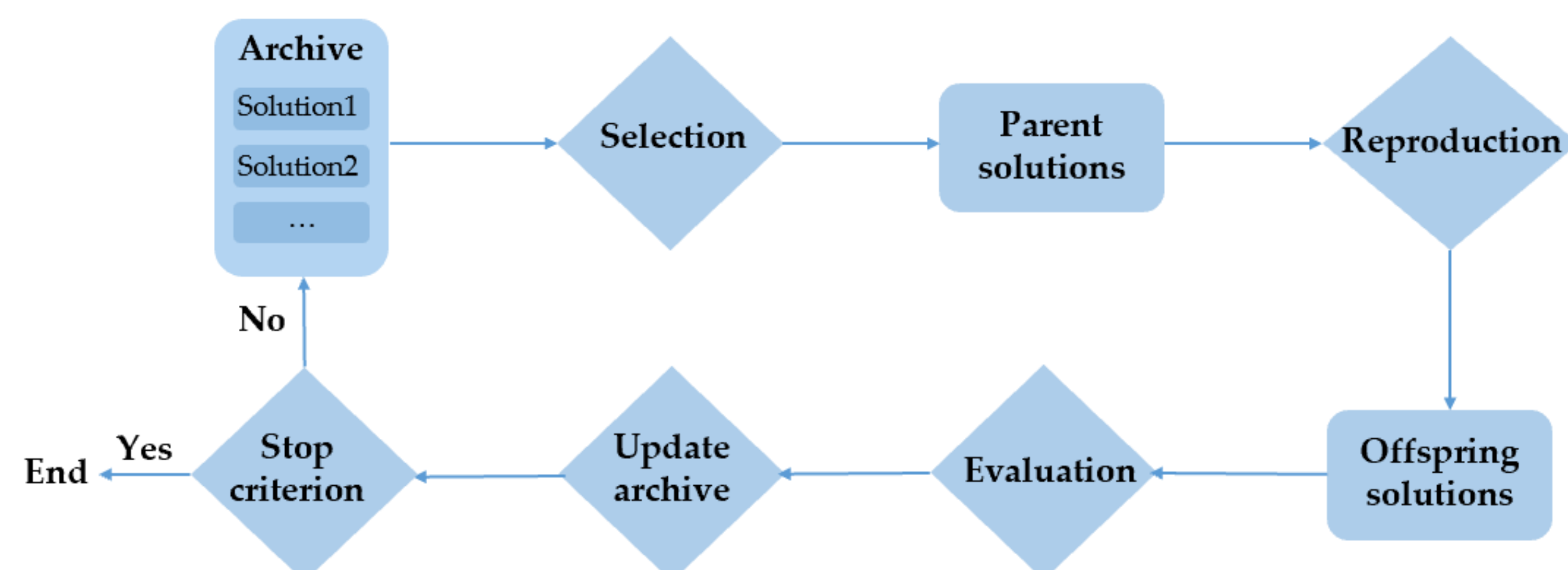
### Reinforcement Learning (RL)

- general RL methods obtain a single policy
- some complex scenarios need a set of diverse policies, which lead to
  - better exploration
  - faster few-shot adaption
  - greater robustness

How to efficiently obtain *a set of high-quality policies with diverse behaviors* is a challenging problem in RL

### Quality-Diversity (QD) algorithms

- a specific type of Evolutionary Algorithms (EAs)
- aims to return a set of high-quality solutions with diverse behaviors

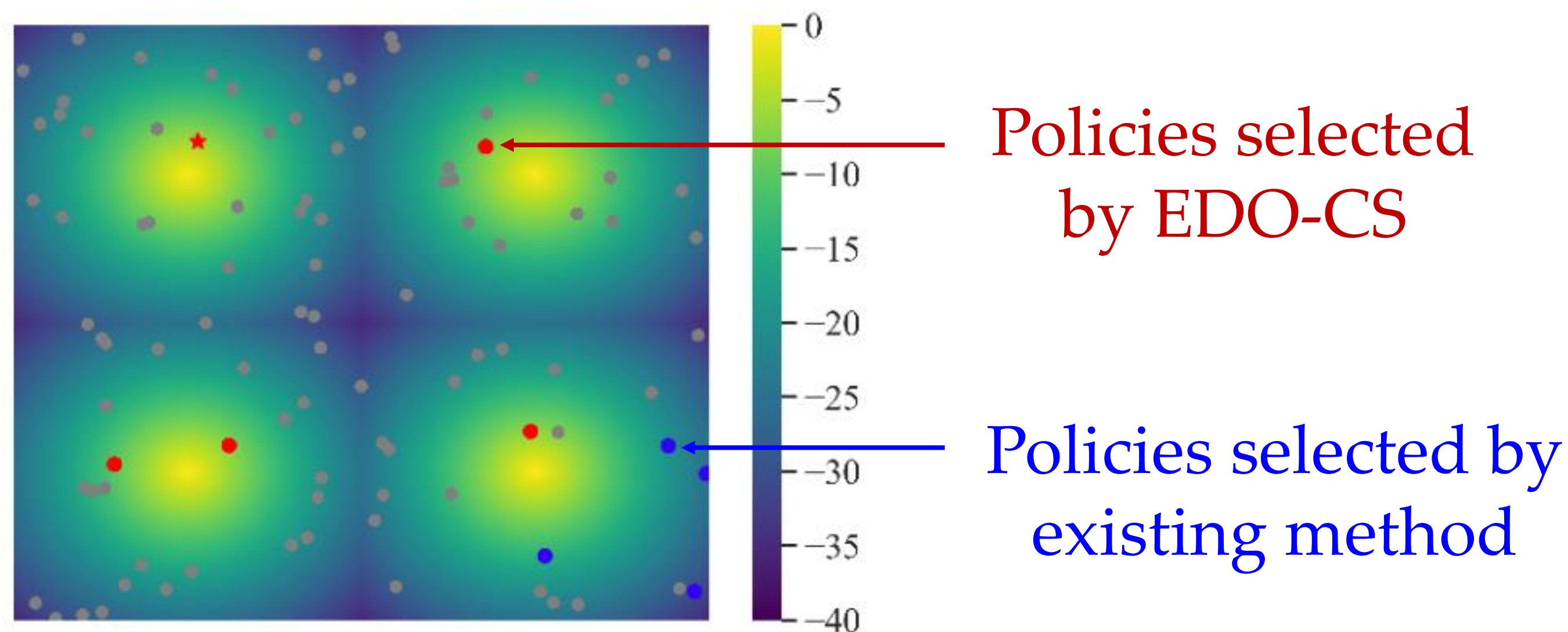


Existing efforts to apply QD algorithms into RL lack an *efficient selection mechanism*

## EDO-CS Method

### Clustering-based selection mechanism

- clusters the policies in the archive based on their behaviors
- selects a high-quality policy from each cluster



### Self-adjusting reproduction mechanism

- the objective function to be maximized

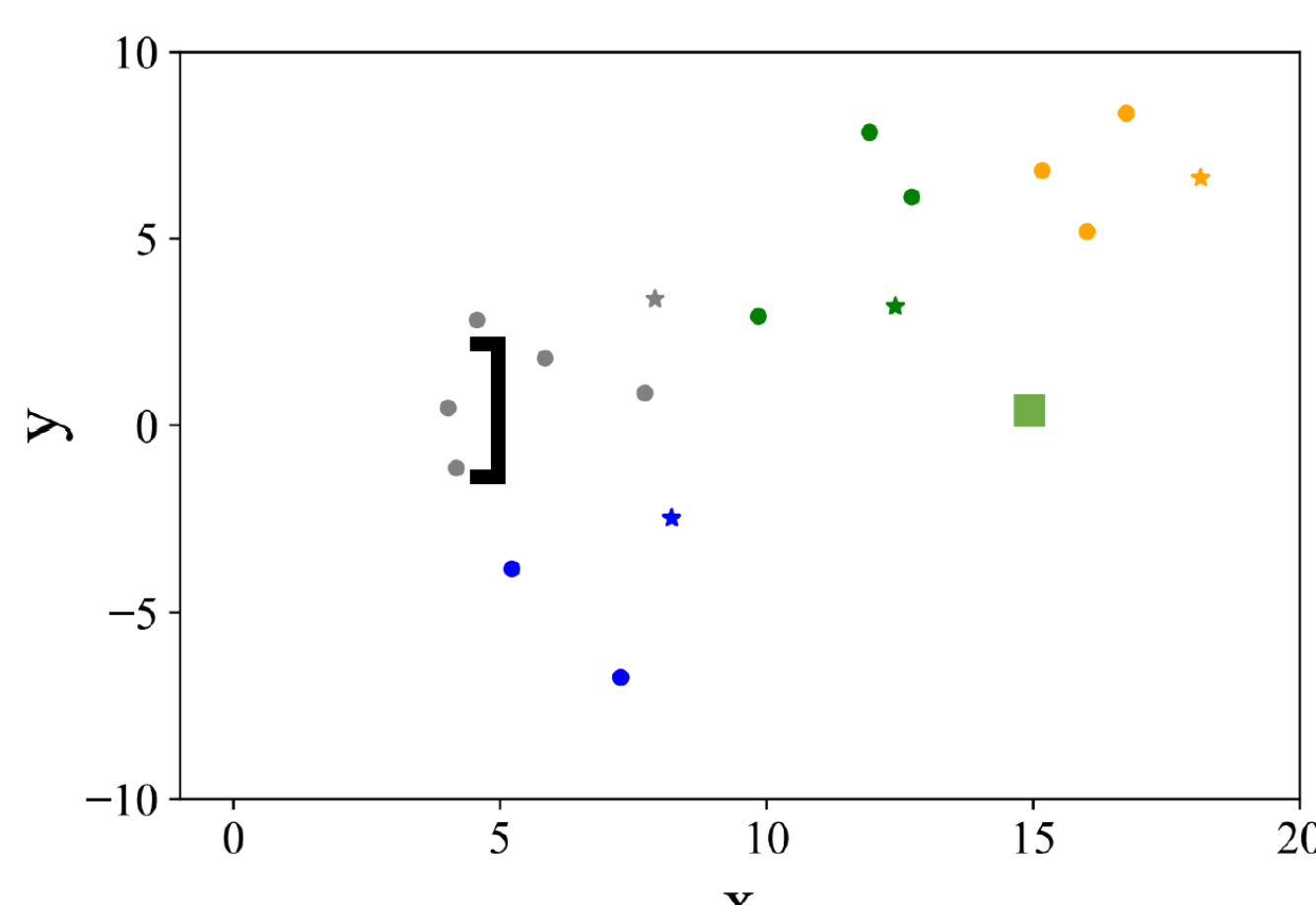
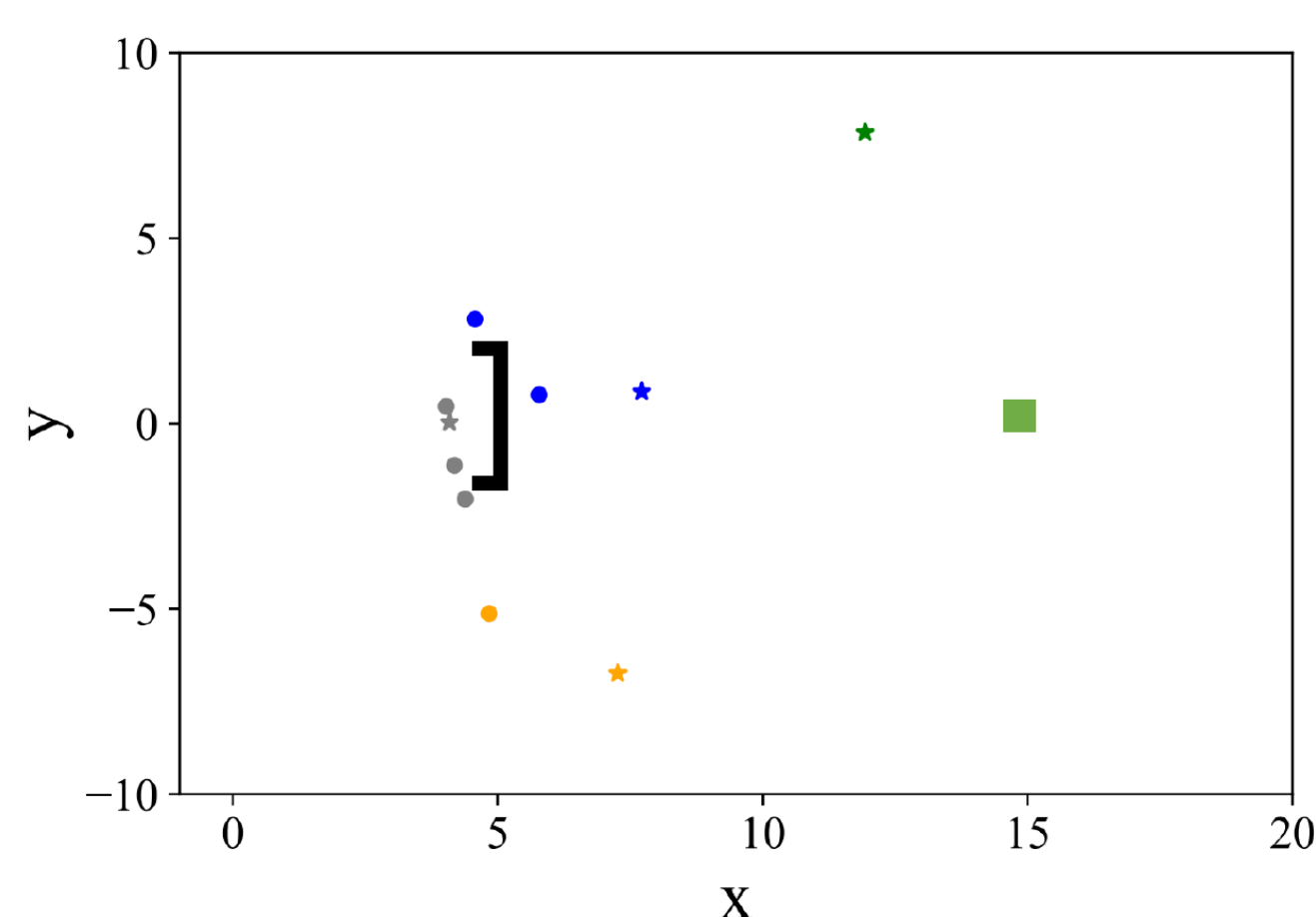
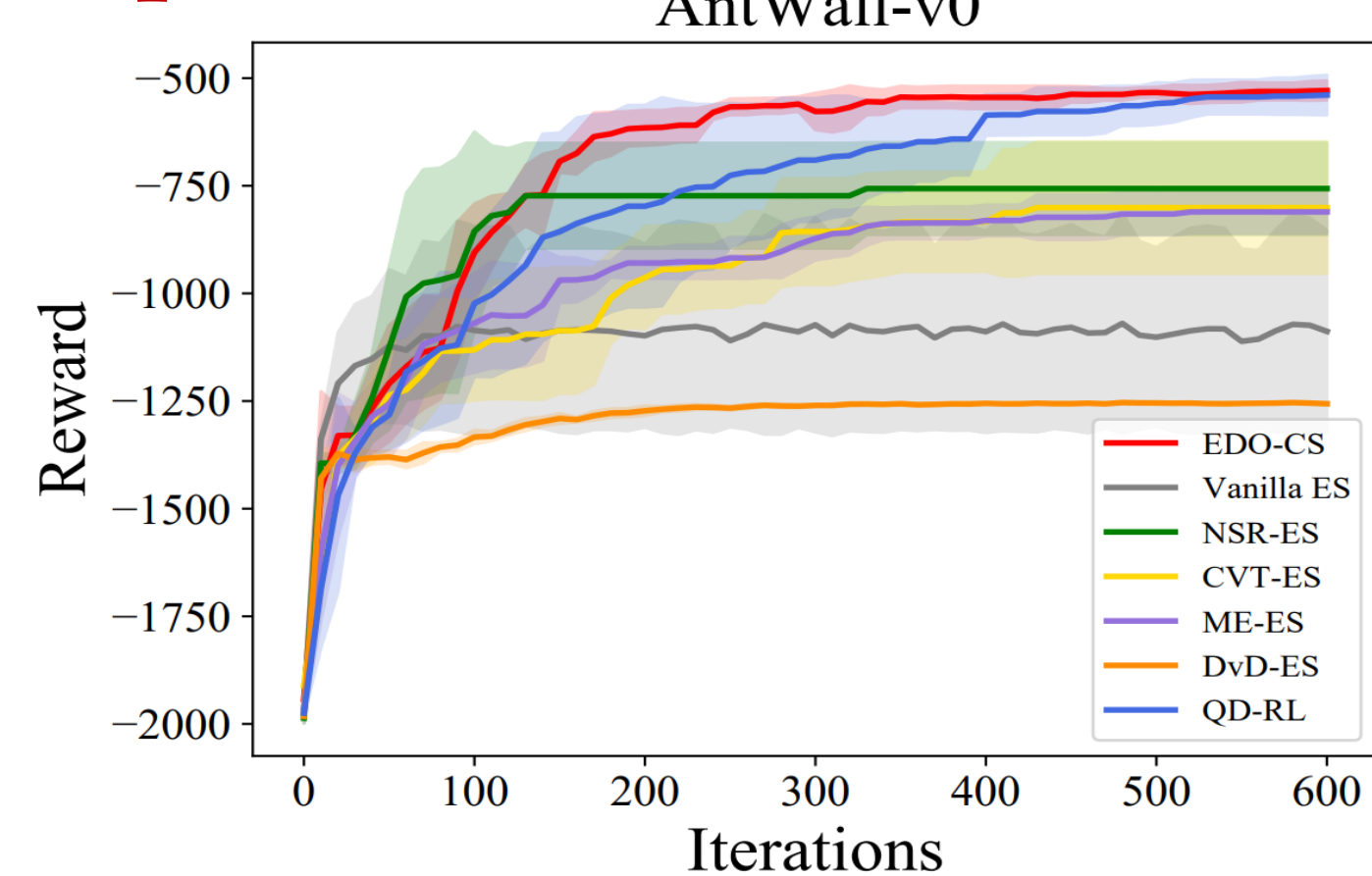
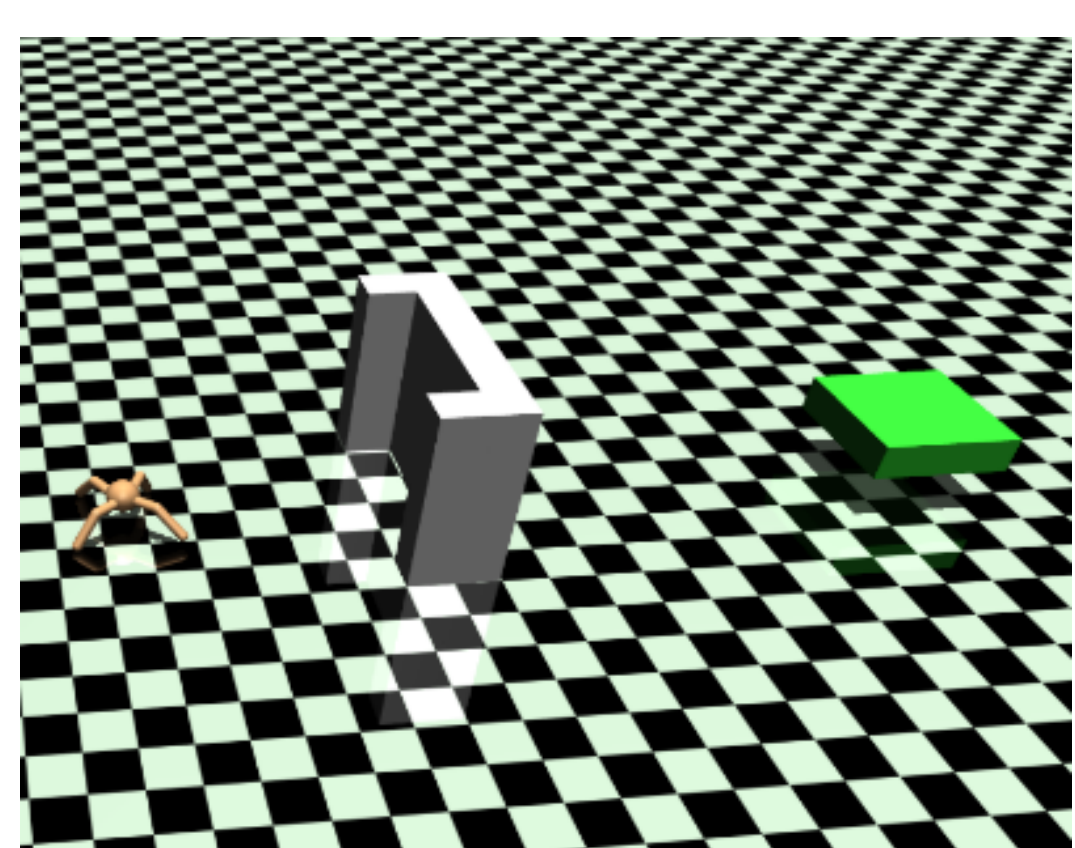
$$J(\theta) = (1 - \lambda)E[R(\tau)] + \lambda Div(\theta)$$

The weight  $\lambda$  controls the trade-off between *exploitation and exploration*, we use multi-armed bandit to self-adjust it

Method	Selection	Reproduction
Vanilla ES	The only parent solution	Quality
NSR-ES	Probabilistic selection	Quality and diversity
CVT-ES	Uniform selection	Quality and diversity
ME-ES	Biased selection	Quality or diversity
DvD-ES	All parent solutions	Quality and diversity
QD-RL	Pareto-based selection	Quality or diversity
<b>EDO-CS</b>	<b>Clustering-based selection</b>	<b>Quality and diversity</b>

## Experiment

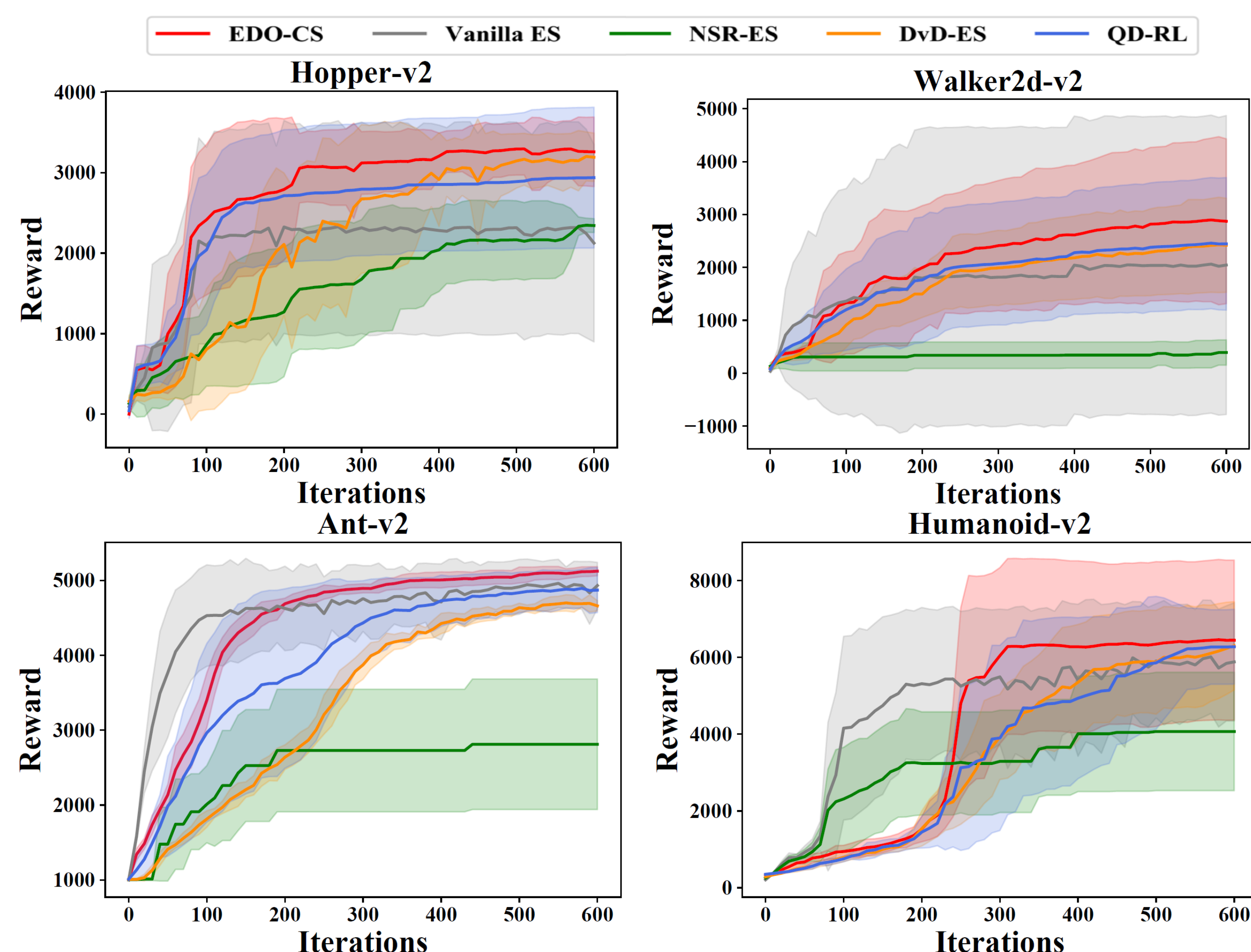
### Task with deceptive rewards



### Multi-modal task

Environment	EDO-CS	QD-RL	ME-ES	DvD-ES	CVT-ES	NSR-ES	Vanilla ES
HalfCheetahFwd	<b>4284</b>	2930	2700	-3419	3219	1346	-5543
HalfCheetahBwd	<b>6548</b>	6013	5953	6353	4672	5366	3911
AntFwd	<b>4617</b>	4291	4316	4507	3856	1737	1911
AntBwd	<b>4697</b>	4164	4123	3498	2958	3961	-851
Performance Ranking	<b>1</b>	3	3.5	3.75	4.75	5.25	6.75

### Single-modal task



EDO-CS shows *superior performance* on various continuous control tasks