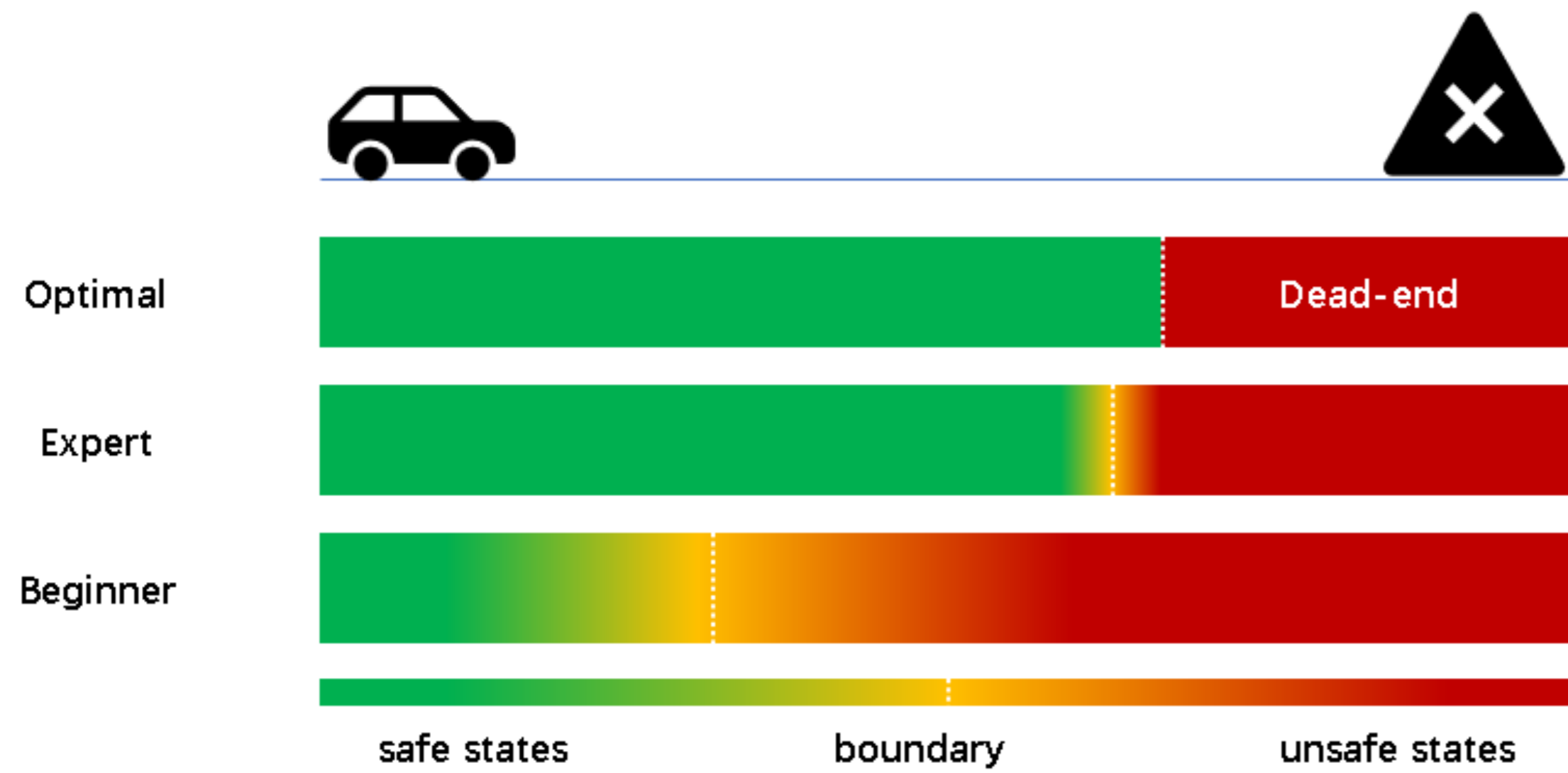


## Motivation

- Existing safe-based methods tend to adopt overly conservative policies.
- The reason is that they assess the safety of states and actions use the task policy.
- The epistemic bias cannot be corrected without violating safety.

## Dead-ends Recognition



Measuring the safety of states **using the optimal safe policy**, which is trained to **minimize the cost and ignore the environmental reward**,

$$V_c^\pi(s) = \mathbb{E}_{\tau \sim \pi, P} \left[ \sum_{t=0}^{\infty} \gamma_{safe}^t \mathcal{C}(s_t, a_t, s_{t+1}) \mid s_0 = s \right]$$

It's possible to **partition the dead-ends** in a deterministic Markov setting **using the state cost function corresponding to the optimal safety policy**.

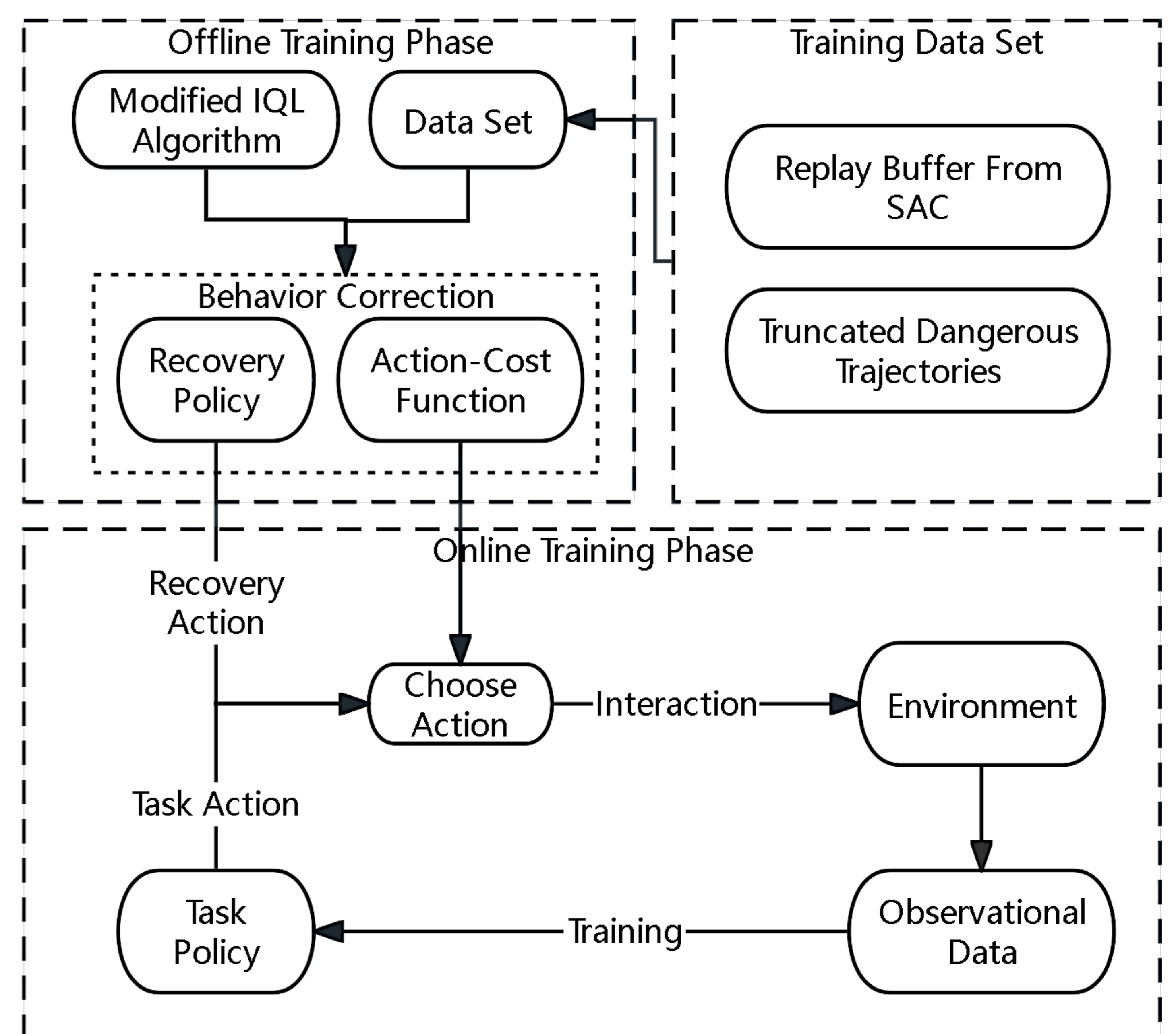
$$S_{dead} = \{s \mid s \in \mathcal{S} \text{ and } \gamma_{safe}^{H-1} \leq V_c^*(s) \leq 1\}$$

## Dead-ends Avoidance

$$A(s)_{safe} = \{a \mid a \in \mathcal{A} \text{ and } Q_c^*(s, a) < \gamma_{safe}^H\}$$

$$a_t = \begin{cases} a^{\pi_{task}}, & Q_c^{\bar{\pi}}(s_t, a^{\pi_{task}}) < \epsilon_{safe} \\ a^{\pi_{rec}}, & \text{otherwise} \end{cases}$$

## Algorithm



*Theorem 1:*  $\Pi_c^{\pi_{com}}$  and  $\Pi_c^{\pi_{rec}}$  are the accessible space for  $\pi$  to explore safely in RRL and DEARRL respectively, we have  $\Pi_c^{\pi_{com}} \subseteq \Pi_c^{\pi_{rec}}$ .

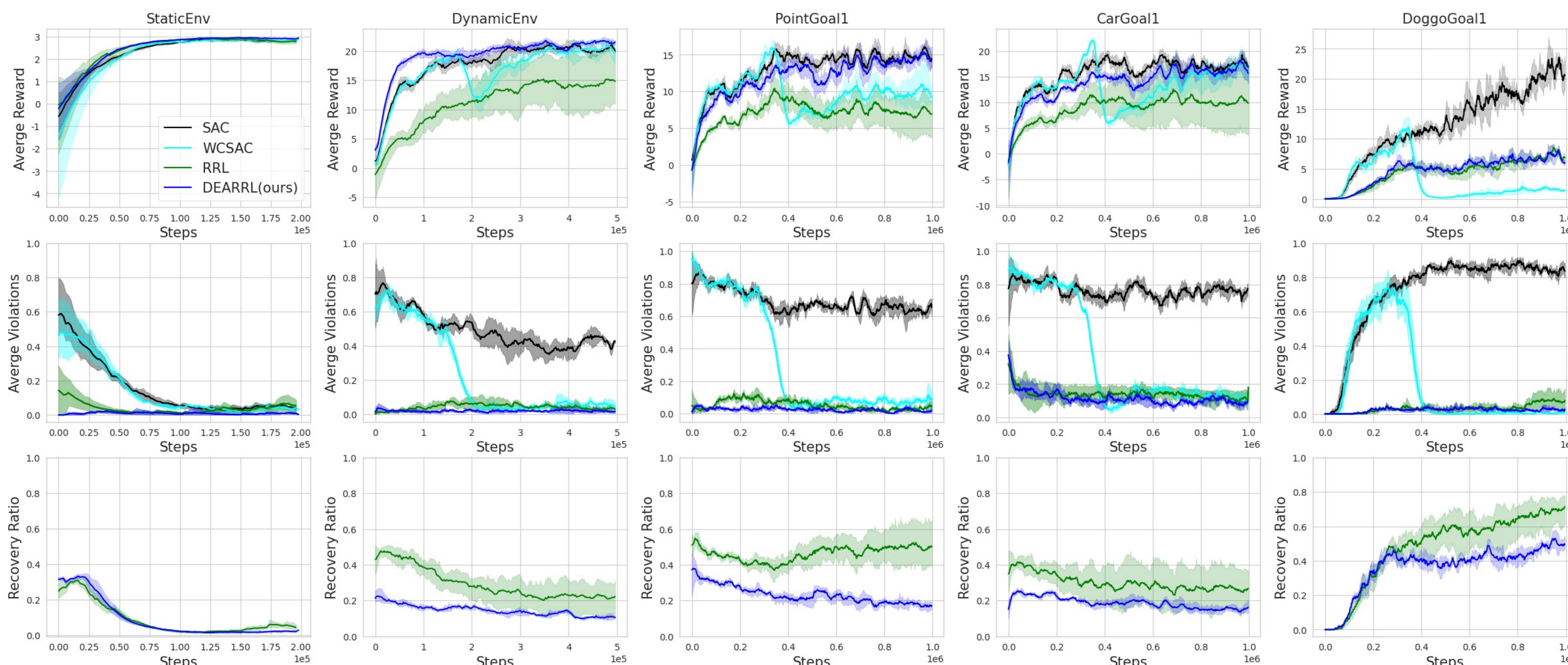
TABLE I  
RESULTS OF TRAINING AND TESTING EACH METHOD INDIVIDUALLY

Environments	SAC			WCSAC			RRL			DEARRL(ours)			IQL	
	ACR	AVR	TV	ACR	AVR	TV	ACR	AVR	TV	ACR	AVR	TV	ACR	AVR
StaticEnv	2.754	0.043	1642	<b>2.936</b>	<b>0.008</b>	1453	2.813	0.032	<u>392</u>	2.882	<u>0.023</u>	<b>119</b>	2.630	0.095
DynamicEnv	17.434	0.519	3330	<u>19.004</u>	0.069	1451	14.230	<b>0.006</b>	<u>112</u>	<b>20.784</b>	<u>0.018</u>	<b>86</b>	18.044	0.565
PointGoal1	12.609	0.753	12375	<u>7.811</u>	<u>0.055</u>	5739	7.203	0.094	<u>476</u>	<b>14.079</b>	<b>0.035</b>	<b>236</b>	<u>12.740</u>	0.76
CarGoal1	16.727	0.73	14850	<b>18.672</b>	<b>0.035</b>	6199	9.647	<u>0.057</u>	<u>1109</u>	<u>15.430</u>	0.091	<b>786</b>	11.357	0.94
DoggoGoal1	16.956	0.908	31450	1.055	<b>0.013</b>	10825	5.924	0.192	<b>623</b>	7.002	<u>0.048</u>	<u>706</u>	<b>18.889</b>	0.838

\* ACR, AVR, TV respectively stand for Average Cumulative Return, Average Violation Rate, and Total Violations in training.

\* The best and second-best results have been marked in bold and underlined, respectively.

## Results



Open-sourced!

<https://github.com/tiev-tongji/dea-rrl>

