

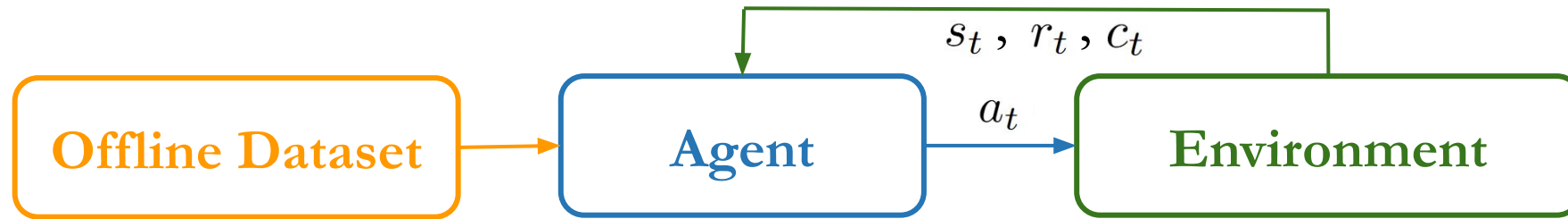
Constraint-Conditioned Actor-Critic for Offline Safe Reinforcement Learning

Zijian Guo¹, Weichao Zhou², Shengao Wang¹, Wenchao Li^{1,2}

¹Division of Systems Engineering, Boston University

²Department of Electrical and Computer Engineering, Boston University
{zjguo, zwc662, wsashawn, wenchao}@bu.edu

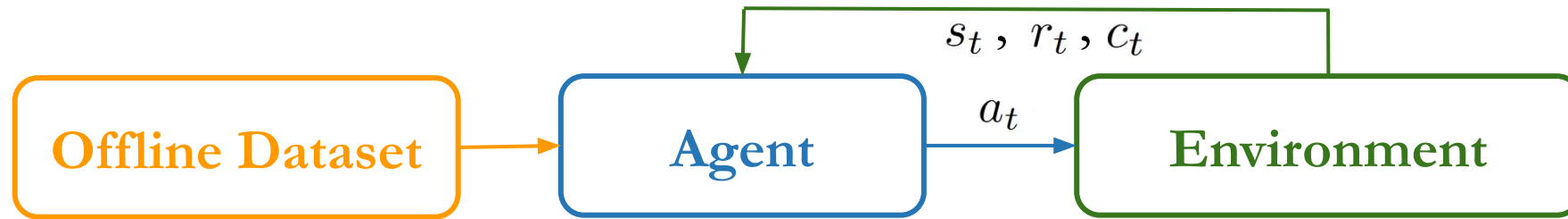
Offline Safe RL



- Constrained MDP: $(\mathcal{S}, \mathcal{A}, P, r, c, \gamma, \mu)$
- Offline dataset: $\mathcal{D} = \{\tau_i\}_{i=1}^N$ with $\tau = \{s_t, a_t, c_t, r_t\}_{t=1}^T$
- Policy that maximizes reward while keeping cost below certain threshold

$$\max_{\pi} \mathbb{E}_{\tau \sim \pi}[R(\tau)] \text{ s.t. } \mathbb{E}_{\tau \sim \pi}[C(\tau)] \leq \epsilon$$

Offline Safe RL



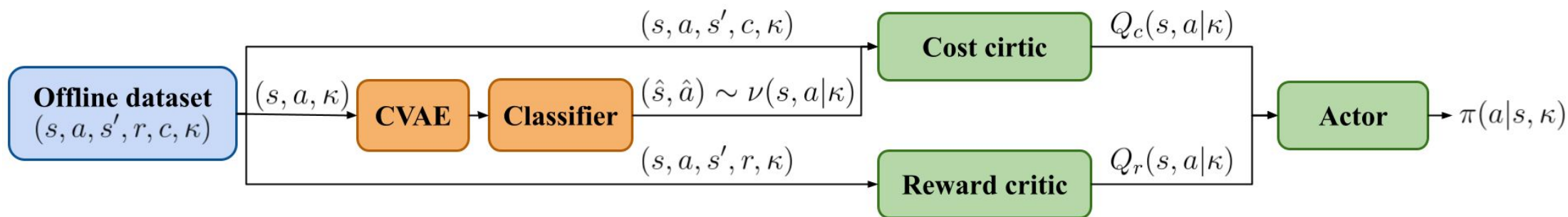
- Constrained MDP: $(\mathcal{S}, \mathcal{A}, P, r, c, \gamma, \mu)$
- Offline dataset: $\mathcal{D} = \{\tau_i\}_{i=1}^N$ with $\tau = \{s_t, a_t, c_t, r_t\}_{t=1}^T$
- Policy that maximizes reward while keeping cost below certain threshold

$$\max_{\pi} \mathbb{E}_{\tau \sim \pi}[R(\tau)] \text{ s.t. } \mathbb{E}_{\tau \sim \pi}[C(\tau)] \leq \epsilon$$

- However, the threshold ϵ is typically **a constant and fixed for training/evaluation**
- **Varying constraint thresholds**
 - Different (or even unseen) thresholds for training/evaluation
 - Time-variant thresholds ϵ_t

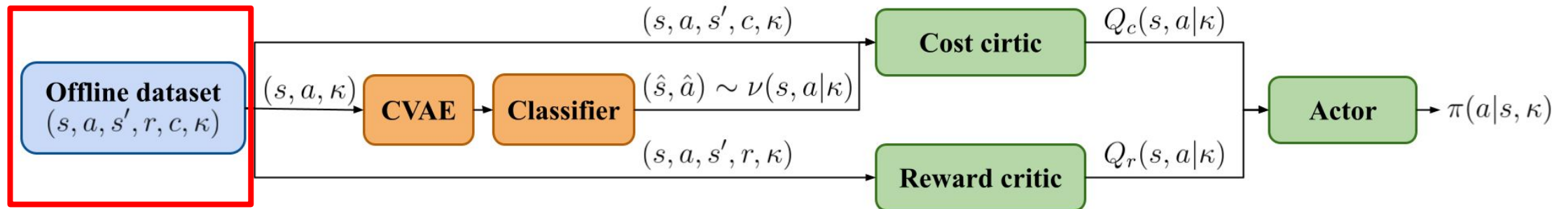
Overview

- **Goal:** learn safe, high-reward, and adaptable policies
- **Key idea:** model the distribution of states and actions and uncover the relationships between behaviors and constraint thresholds from the offline dataset
- **Constraint-conditioned actor-critic (CCAC)**



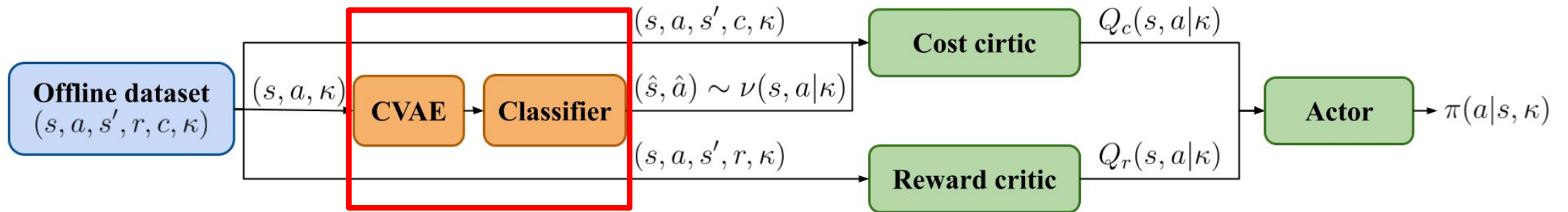
Overview of CCAC

Varying Constraint Thresholds



- Define the varying constraint threshold as **cost budget** κ_t
 - Set the cumulative cost of each trajectory as its initial cost budget: $\kappa_1 = \sum_{t=1}^T c_t$
 - Updated based on incurred cost: $\kappa_{t+1} = \kappa_t - c_t$
- Given a cost budget κ , a state-action pair (s, a) is considered **out-of-distribution (OOD)** if taking the action a at the state s will lead to exceeding the cost budget κ

OOD Distribution Modeling



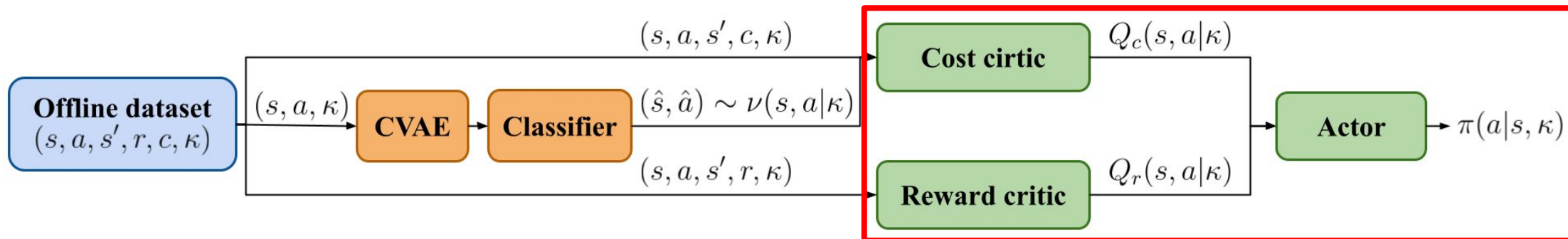
- Constraint-conditioned variational autoencoder (CVAE):

$$\max_{p,q} \mathbb{E}_{s,a,\kappa \sim \mathcal{D}} \left[\mathbb{E}_{z \sim q(z|s,a,\kappa)} [p(s, a|z, \kappa)] - \beta D_{KL}(q(z|s, a, \kappa) || p(z|\kappa)) \right]$$

where $q(z|s, a, \kappa)$ is the encoder and $p(s, a|z, \kappa)$ is the decoder

- Classifier: $\min_h \mathbb{E}_{s,a,\kappa \sim \mathcal{D}} \left[-y \log(h(s, a|\kappa)) - (1 - y) \log(1 - h(s, a|\kappa)) \right]$
- The OOD distribution $\nu(s, a|\kappa)$ is proportional to $\mathbb{E}_{z \sim p(z|\kappa)} [p(s, a|z, \kappa) h(s, a|\kappa)]$**

Constraint-Conditioned Actor-Critic



- Constraint-conditioned cost critic (**overestimation of OOD state-action pairs**):

$$\min_{Q_c} \mathbb{E}_{s,a,\kappa \sim \mathcal{D}} [(Q_c(s, a | \kappa) - \mathcal{T}^\pi Q_c(s, a | \kappa))^2] , \text{ s.t. } \mathbb{E}_{s,a \sim \nu, \kappa \sim \mathcal{D}} [Q_c(s, a | \kappa)] \geq \epsilon$$

- Constraint-conditioned reward critic:

$$\min_{Q_r} \mathbb{E}_{s,a,\kappa \sim \mathcal{D}} [(Q_r(s, a | \kappa) - \mathcal{T}^\pi Q_r(s, a | \kappa))^2]$$

- Constraint-conditioned actor:

$$\max_{\pi} \mathbb{E}_{s,\kappa \sim \mathcal{D}, a \sim \pi(a | s, \kappa)} [Q_r(s, a | \kappa)] , \text{ s.t. } \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi(a | s, \kappa)} [Q_c(s, a | \kappa)] \leq \kappa, \forall \kappa \in \mathcal{D}$$

Experimental Settings

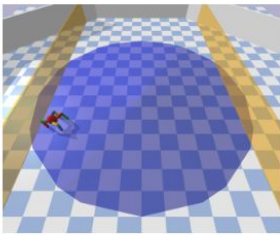
- **Environments** (BulletSafetyGym[1], SafetyGymnasium[2])

- **Three tasks: Run, Circle, and Velocity.**

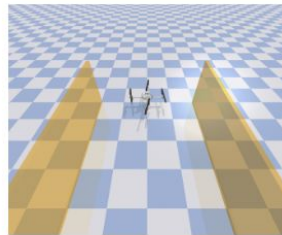
Run/Velocity: run as fast as it can within the speed limit between two boundaries

Circle: run in a circle but are constrained within a safe region

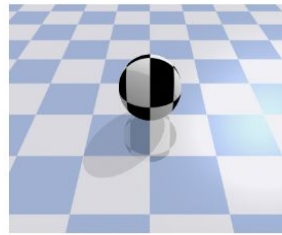
- **Multiple robots:** Ant, Ball, Car, Drone, Hopper, HalfCheetah



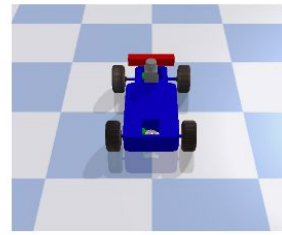
Circle task



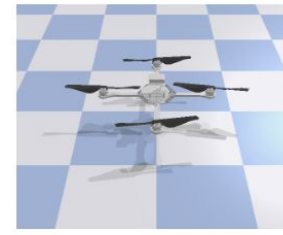
Run task



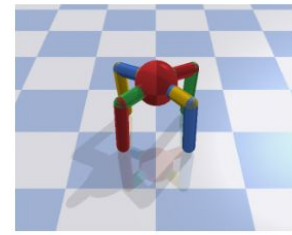
Ball



Car



Drone



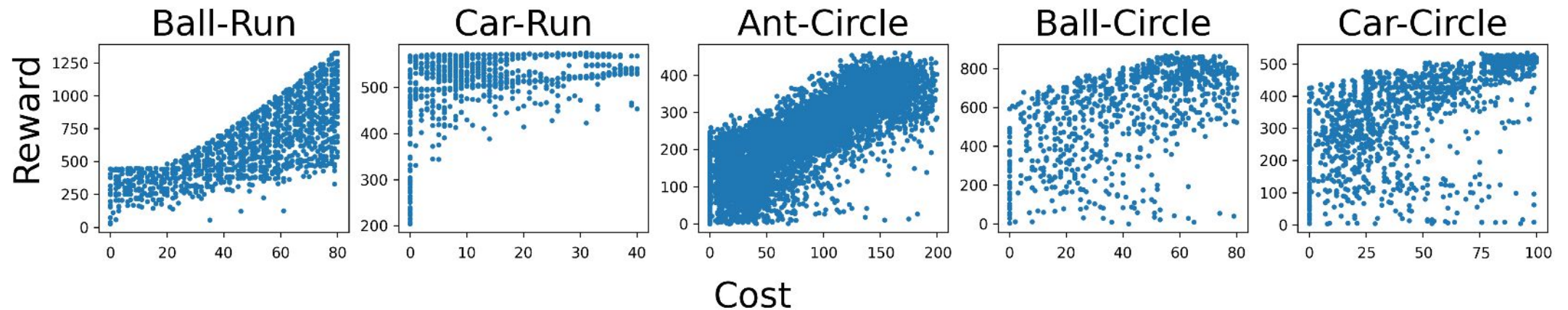
Ant

[1] Gronauer, Sven. "Bullet-safety-gym: A framework for constrained reinforcement learning." (2022).

[2] Ji, Jiaming, et al. "Safety gymnasium: A unified safe reinforcement learning benchmark." *Advances in Neural Information Processing Systems* 36 (2023).

Experimental Settings

- **Offline datasets: DSRL[3]**
 - Reward v.s. cost plots, each dot represents a trajectory



[3] Liu, Zuxin, et al. "Datasets and benchmarks for offline safe reinforcement learning." *arXiv preprint arXiv:2306.09303* (2023).

Results

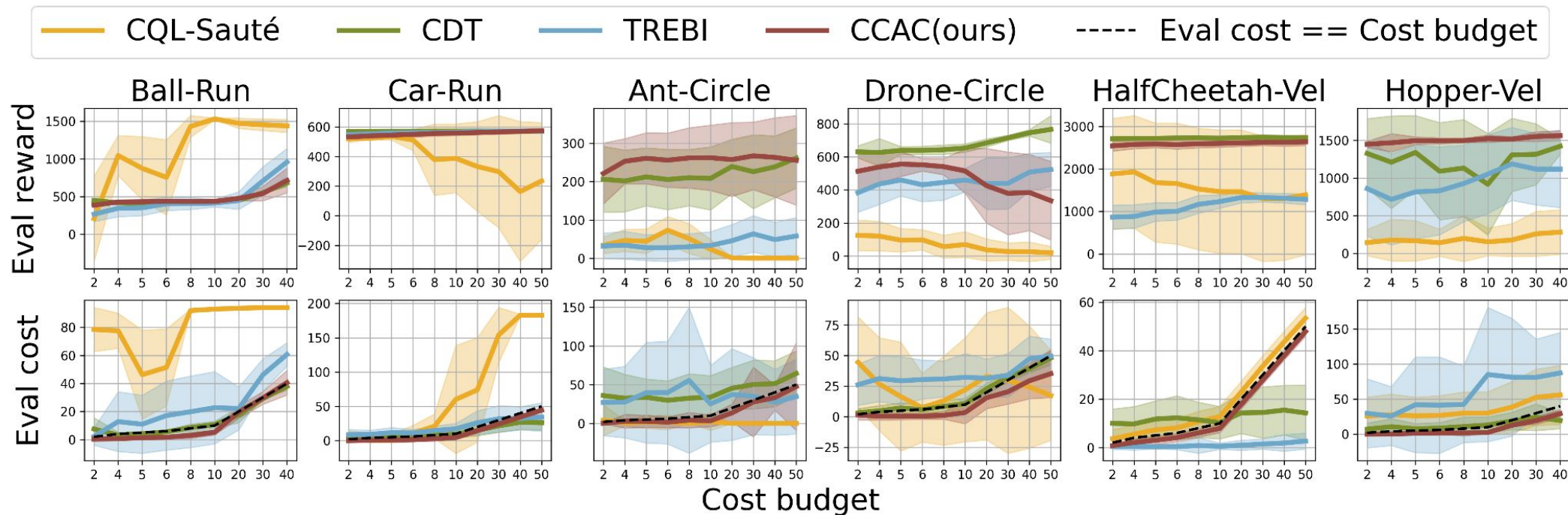
- CCAC can learn **safe and high-reward policies** from offline datasets

Tasks	Metric	CQL-Sauté	BCQ-Lag	BEAR-Lag	CPQ	COptiDICE	VOCE	CDT	TREBI	FISOR	CCAC(ours)
Run	Reward \uparrow	0.55 \pm 0.44	1.51 \pm 0.91	1.88 \pm 1.03	0.95 \pm 0.03	1.38 \pm 0.5	1.65 \pm 1.07	0.99 \pm 0.02	0.82 \pm 0.19	0.74\pm0.08	0.96\pm0.03
	Cost \downarrow	5.23 \pm 6.26	19.85 \pm 12.76	23.56 \pm 7.99	1.92 \pm 2.66	6.27 \pm 6.29	9.35 \pm 8.95	1.34 \pm 0.74	2.12 \pm 2.36	0.54\pm1.90	0.23\pm0.27
Circle	Reward \uparrow	0.4 \pm 0.33	1.11 \pm 0.31	1.13 \pm 0.21	0.64 \pm 0.42	0.63 \pm 0.25	0.06 \pm 0.11	0.91 \pm 0.17	0.43 \pm 0.25	0.27\pm0.18	0.79\pm0.24
	Cost \downarrow	6.88 \pm 9.15	15.61 \pm 6.1	13.19 \pm 6.18	3.73 \pm 6.43	9.82 \pm 10.57	10.05 \pm 15.53	2.58 \pm 2.49	3.53 \pm 8.26	0.17\pm0.76	0.17\pm0.79
Velocity	Reward \uparrow	0.44 \pm 0.43	0.78 \pm 0.37	0.18 \pm 0.8	0.41 \pm 1.15	0.65 \pm 0.36	-0.41\pm0.44	0.88 \pm 0.24	0.45 \pm 0.24	0.47\pm0.26	0.86\pm0.2
	Cost \downarrow	2.08 \pm 3.23	27.26 \pm 29.75	20.1 \pm 25.76	35.18 \pm 44.66	6.94 \pm 7.22	0.0\pm0.0	2.91 \pm 2.87	1.07 \pm 2.52	0.08\pm0.27	0.38\pm0.2
Average	Reward \uparrow	0.45 \pm 0.4	1.09 \pm 0.58	0.98 \pm 0.93	0.63 \pm 0.75	0.8 \pm 0.47	0.26 \pm 0.96	0.91 \pm 0.19	0.51 \pm 0.28	0.44\pm0.27	0.85\pm0.21
	Cost \downarrow	4.91 \pm 7.33	20.44 \pm 19.33	17.8 \pm 16.46	13.81 \pm 30.22	8.07 \pm 8.85	6.54 \pm 12.1	2.41 \pm 2.45	2.41 \pm 6.0	0.22\pm1.05	0.25\pm0.56

Bold: safe agents with (normalized) cost smaller than 1. **Blue:** safe agents with the highest reward

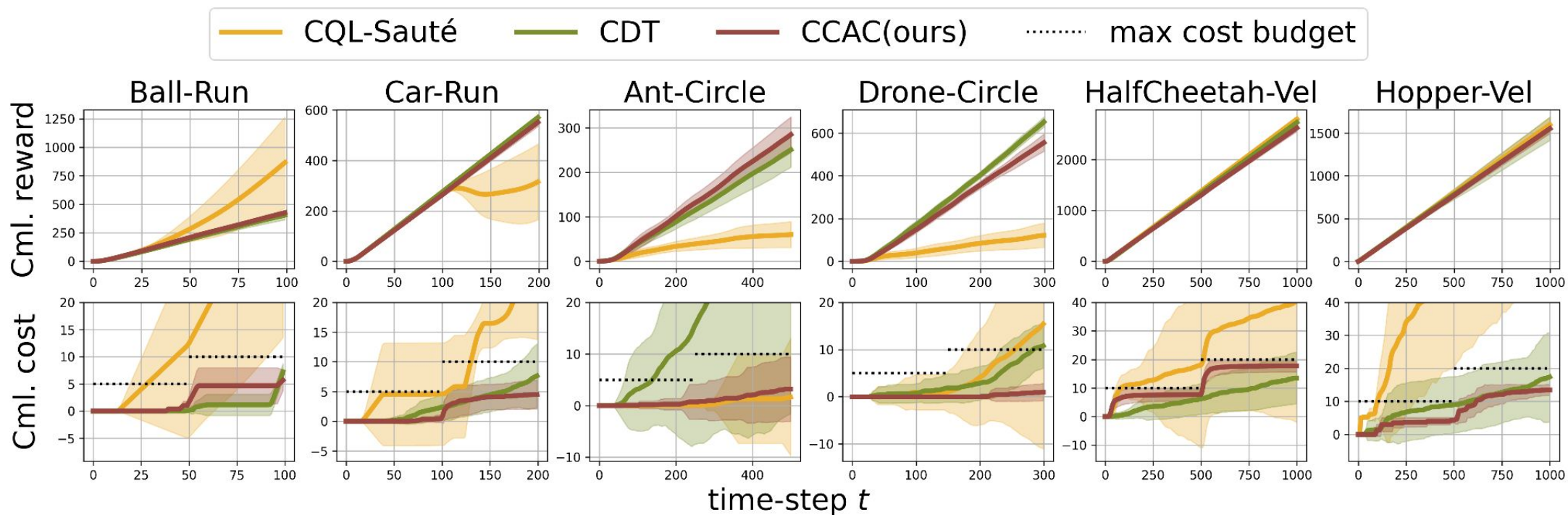
Results

- CCAC can achieve **zero-shot adaptation to different (even unseen) cost thresholds**



Results

- CCAC can achieve **zero-shot adaptation to dynamic cost thresholds**



Summary

- A novel OSRL method CCAC:
 - Handle varying constraint thresholds
 - Model the distribution of states and actions based on constraint thresholds
 - Overestimate the cost critics of OOD state-action pairs
- CCAC is able to learn **safe**, **high-reward** and **adaptable** policies
- Code is available at <https://github.com/BU-DEPEND-Lab/CCAC>