



Computer  
Science



---

# Off-Policy Deep Reinforcement Learning with Analogous Disentangled Exploration

---

Anji Liu<sup>1</sup>, Yitao Liang<sup>1</sup>, Guy Van den Broeck<sup>1</sup>

<sup>1</sup>Computer Science Department, UCLA

# Motivation: Why Reinforcement Learning

---

Is capable of solving **large-scale and complex problems**.

Learning through a trial-and-error process **with little supervision**.



Video games

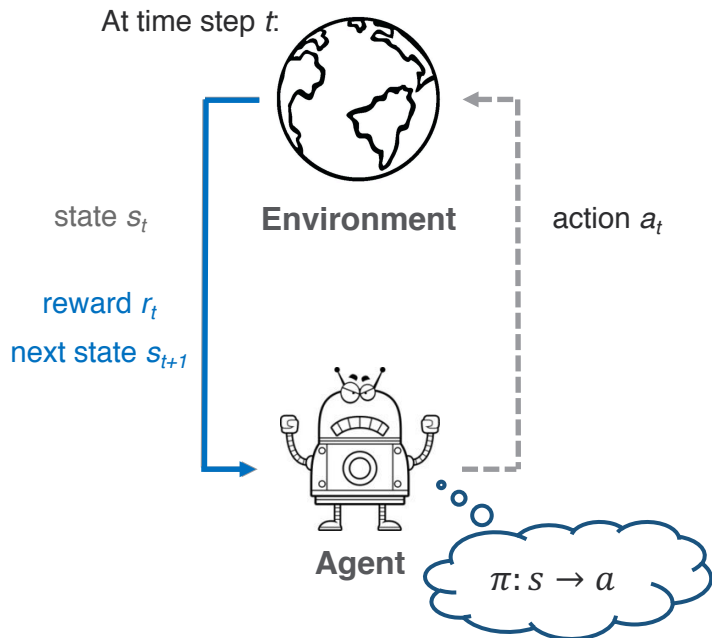


Autonomous driving



Go

# Background: Markov Decision Process



## Policy

Determine which action to take given a state.

## Value function

Measures discounted long-term reward

$$Q(s, a) = \mathbb{E}_{a_t \sim \pi, s_{t+1} \sim \mathcal{P}, r_t \sim \mathcal{R}} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a \right]$$

## Goal of RL

Find a policy  $\pi$  that maximizes the expected long-term reward.

$$\mathbb{E}_{s \sim \rho_0, a \sim \pi} [Q(s, a)]$$

where  $\rho_0$  is the initial state distribution.

# The Exploration-Exploitation Tradeoff

---



The necessity of two policies<sup>1</sup>: **target policy**  $\pi$  and **behavior policy**  $\mu$

<sup>1</sup>For simplicity, sometimes only one policy is explicitly constructed.

# Goals

---

## What we want to achieve?

Leverage the flexibility to **explicitly design two policies** to **better balance the exploration-exploitation tradeoff**.

## How we achieve it?

Analogous Disentangled Actor Critic

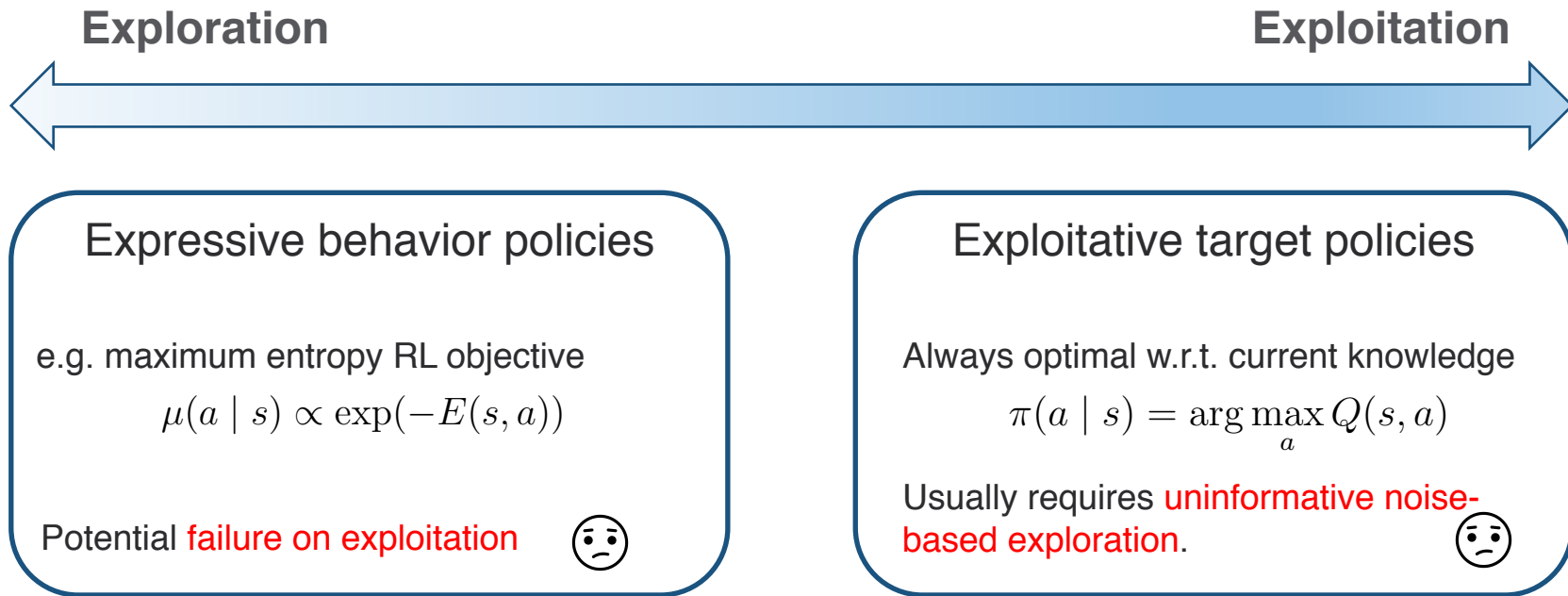


Restricting the disentangled behavior policy (**policy co-training**).



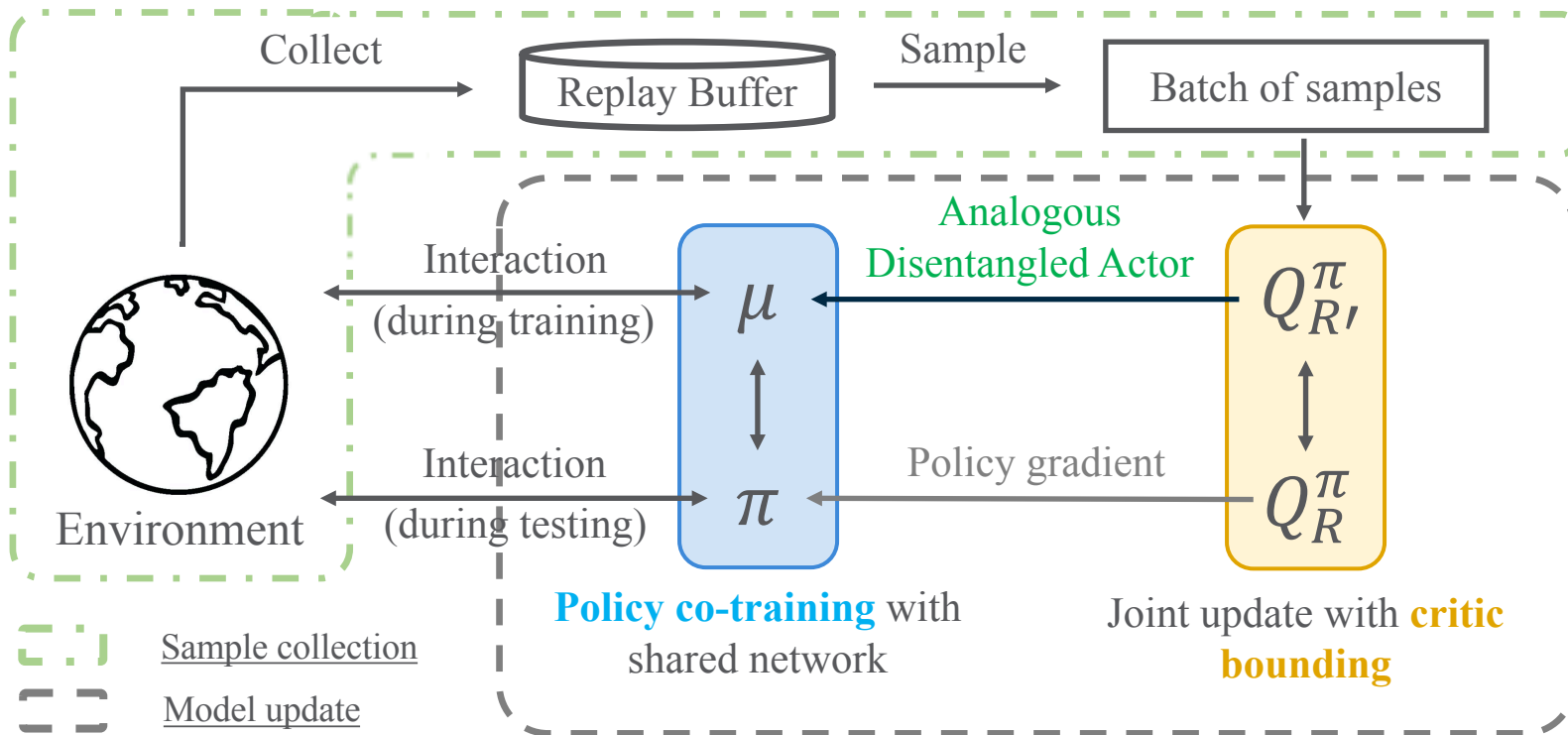
Restricting value updates (**critic bounding**).

# Combating the Exploration-Exploitation Dilemma



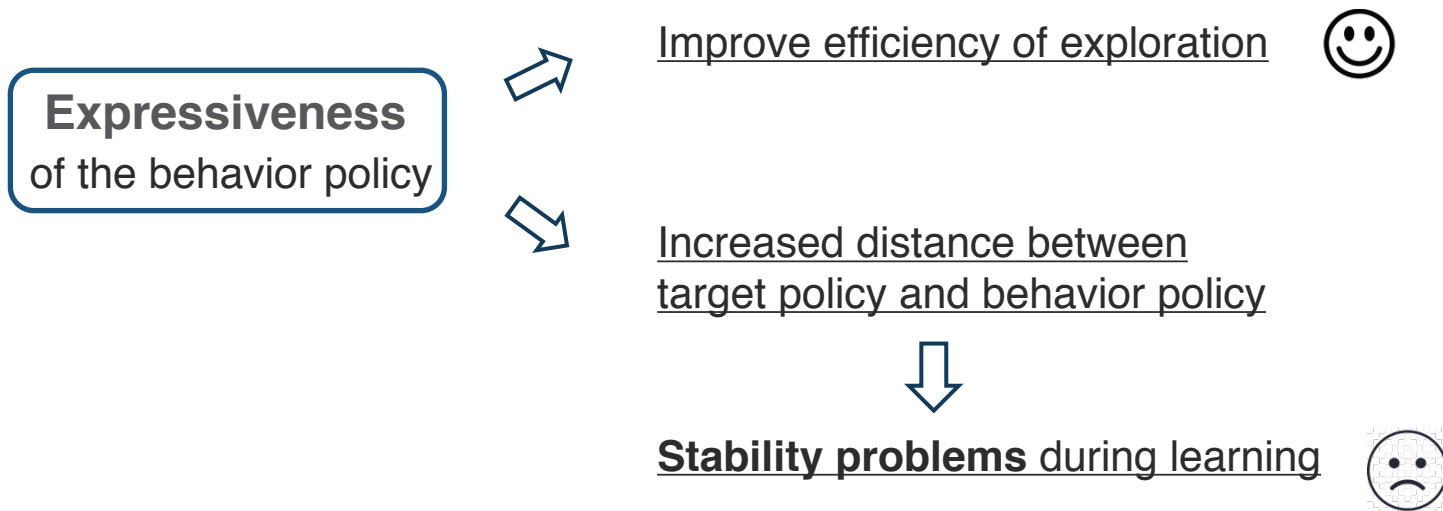
Achieving **better tradeoff** through the flexibility to **design separated policies?**

# Analogous Disentangled Actor-Critic (ADAC)



# Stabilizing Policy Updates by Policy Co-training

## Motivation



Obtain **expressive** while **stable and optimal** behavior policy.



# Stabilizing Policy Updates by Policy Co-training

---

## Key observation

classic policy optimization

$$\pi(a \mid s) = \arg \max_a Q(s, a)$$

Deterministic policy gradient



maximum-entropy policy optimization

$$\mu(a \mid s) \propto \exp(Q(s, a))$$

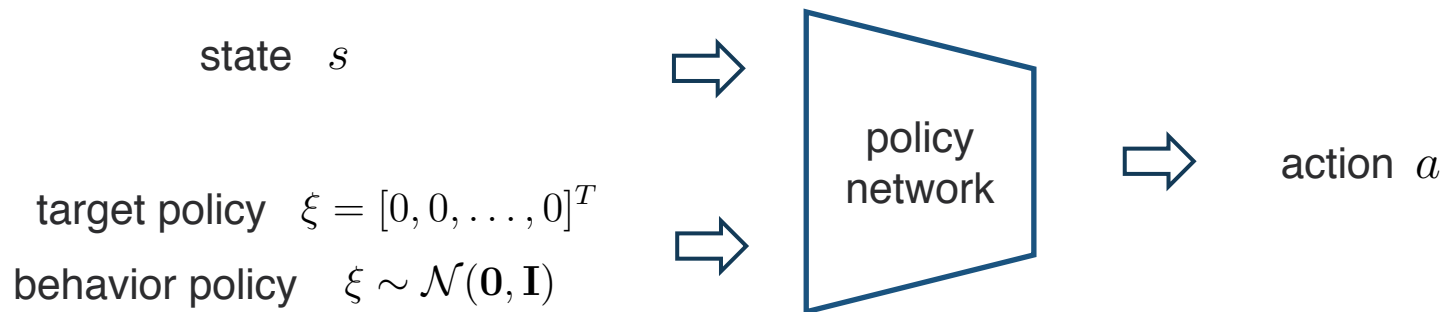
Deterministic policy gradient +  
entropy regularization



**Joint learning** of behavior/target policy in a **shared neural network**.

# Analogous Disentangled Behavior Policy

---



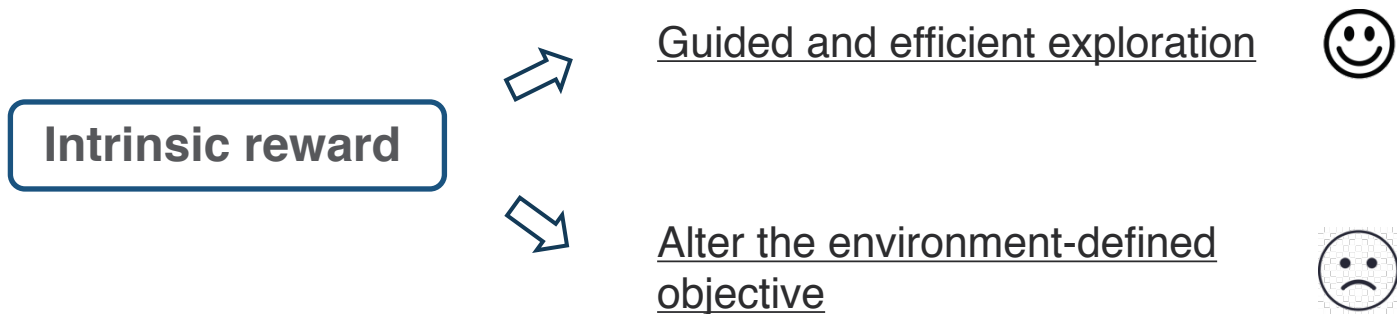
## Benefits

Increase the effectiveness of the behavior policy

- Reduce the **distance** between two policies and stabilize the learning process.
- Allowing **expressive** behavior policy for efficient exploration.

# Incorporating Intrinsic Reward via Critic Bounding

## Motivation



Letting intrinsic reward to **only affect the behavior policy.**

# Incorporating Intrinsic Reward via Critic Bounding

---

$Q_{\mathcal{R}}^{\pi}$  value function w.r.t. the target policy and the environment-defined reward  
→ Adopted to improve the target policy.

$Q_{\mathcal{R}'}^{\pi}$  value function w.r.t. the target policy and the enhanced reward,  
where  $\mathcal{R}' = \mathcal{R} + \mathcal{R}^{\text{in}}$   
→ Adopted to improve the behavior policy.

## Theoretical justification

- Bounded training stability.
- Bounded training effectiveness.

# Summary of Advantages

---

## Policy co-training

Bounds the target policy and the behavior policy to stabilize the training process.

Allowing expressive behavior policy.

## Critic bounding

Incorporate intrinsic reward for effective exploration.

Has no effect on the optimality of the target policy.

# Analysis of Analogous Disentangled Behavior Policy

## Key results

### Behavior policy

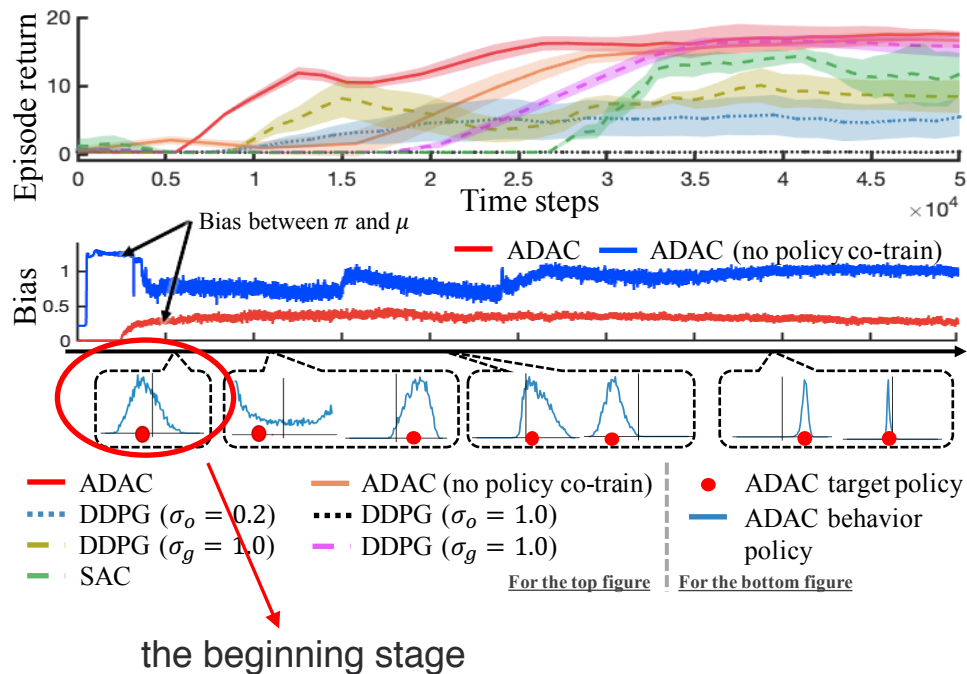
- Act curiously at the beginning.
- Focus on potentially rewarding actions after obtaining preliminary understanding of the environment.

### Target policy

- Remains optimal w.r.t. the current knowledge.

### Target policy & behavior policy

- Bias between them significantly reduced.



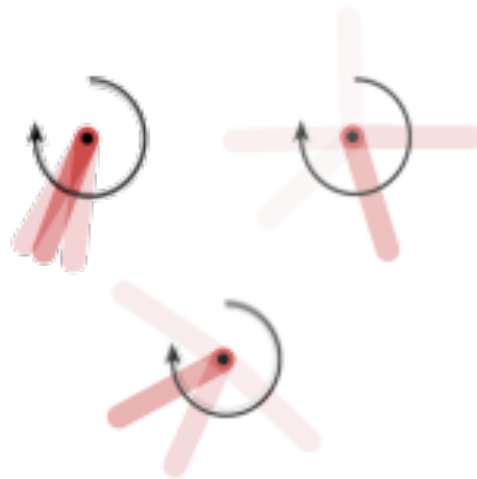
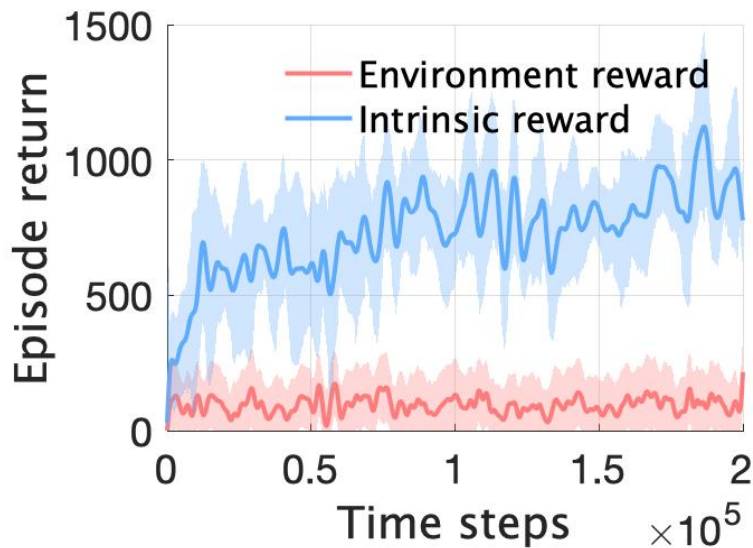
# Comparison with the State of the Art

Out-performs state-of-the-art methods in 10 out of 14 benchmark environments.

- ADAC consistently out-performs its base model, and retains the benefits of improvements developed by the base models.
- Comparison with SAC reveals the benefit brought by the disentangled structure.

Environment	ADAC (TD3)	ADAC (DDPG)	TD3	DDPG	SAC	PPO
RoboschoolAnt	2219±373	838.1*±97.1	<b>2903</b> †±666	450.0±27.9	<b>2726</b> ±652	1280±71
RoboschoolHopper	<b>2299</b> †±333	766.5*±10	<b>2302</b> †±537	543.8±307	2089±657	1229±345
RoboschoolHalfCheetah	1578†±166	<b>1711</b> *±95	607.2±246.2	441.6±120.4	807.0±252.6	1225±184.2
RoboschoolAtlasForwardWalk	<b>234.6</b> †±55.7	186.7*±37.9	190.6±50.1	52.63±26.2	126.0±47.1	107.6±29.4
RoboschoolWalker2d	<b>1769</b> †±452	<b>1564</b> *±651	995.1±146.3	208.7±137.1	1021±263	578.9±231.3
Ant	3353±847	1226*±18	4034†±517	370.5±223	<b>4291</b> ±1498	1401±168
Hopper	<b>3598</b> †±374	374.5*±36.5	2845±609	38.93±0.88	<b>3307</b> ±825	1555±458
HalfCheetah	9392±199	2238*±40	10526†±2367	1009±49	<b>11541</b> ±2989	881.7±10.1
Walker2d	<b>5122</b> †±1314	1291*±42	4630†±778	186.2±33.3	4067±1211	1146±368
InvertedPendulum	<b>1000</b> †±0	<b>1000</b> *±0	<b>1000</b> †±0	<b>1000</b> *±0	<b>1000</b> ±0	98.90±2.08
InvertedDoublePendulum	<b>9359</b> †±0.17	9334*±1.39	7665±566	27.20±2.61	9353±2896	98.90±5.88
BipedalWalker	<b>309.8</b> †±15.6	-52.77*±1.94	288.4†±51.25	-123.90±11.17	<b>307.2</b> ±57.92	266.9±28.52
BipedalWalkerHardcore	<b>-10.76</b> †±27.70	-98.52±3.21	-57.97±21.08	-50.05*±10.27	-127.4±45.2	-105.3±22.2
LunarLanderContinuous	<b>290.0</b> †±50.9	85.67*±23.42	<b>289.7</b> †±54.1	-65.89±96.48	283.3±69.29	59.32±68.44

# Why intrinsic reward sometimes harms RL

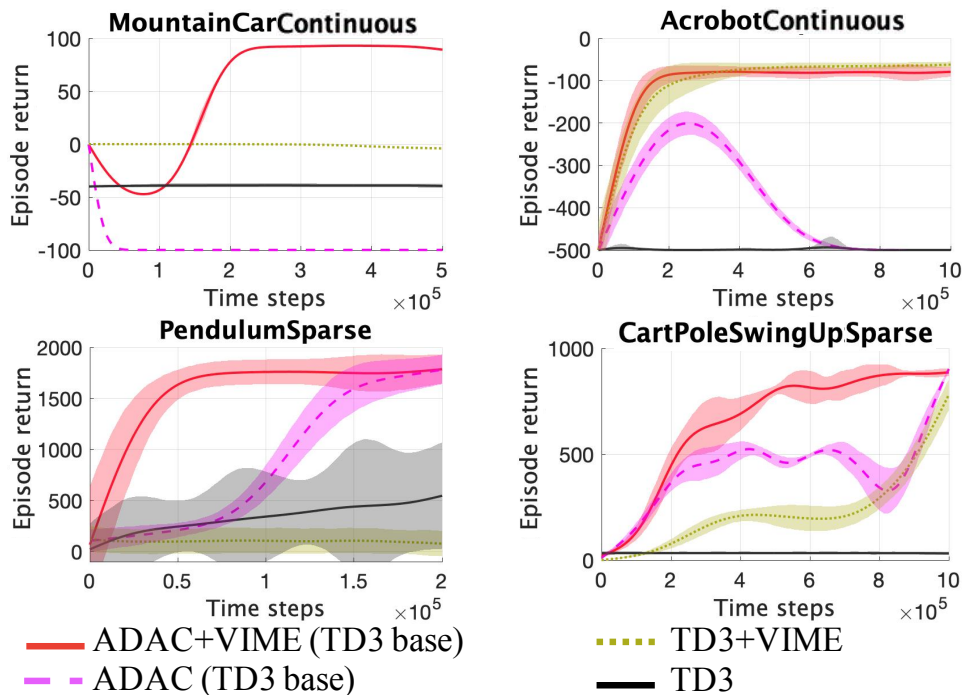


**Environment-defined reward** stays the same while **intrinsic reward** keep growing.



# ADAC with Intrinsic Reward

ADAC out-performs baseline methods on challenging sparse-reward tasks when using intrinsic reward.



# Thank You

---

Open source code: [github.com/UCLA-StarAI/Analogous-Disentangled-Actor-Critic](https://github.com/UCLA-StarAI/Analogous-Disentangled-Actor-Critic)