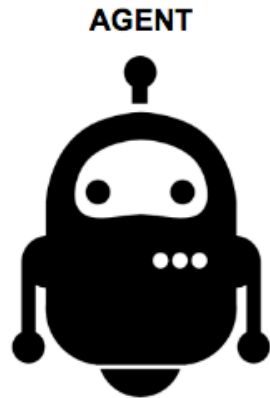


OMPO: A Unified Framework for RL under Policy and Dynamics *Shifts*

Yu Luo, Tianying Ji, Fuchun Sun*, Jianwei Zhang, Huazhe Xu & Xianyuan Zhan



Interacting with environments

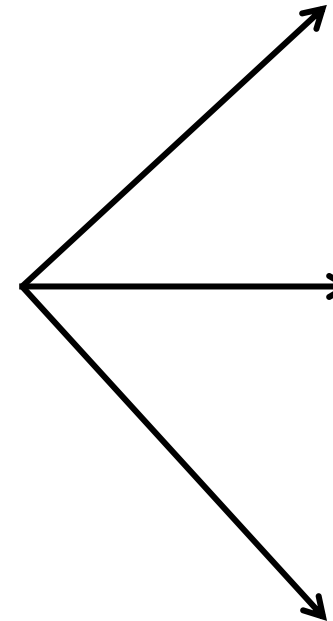


- State $s \in \mathcal{S}$
- Take action $a \in \mathcal{A}$

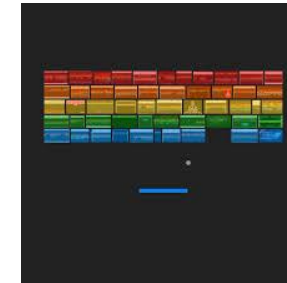
ENVIRONMENT



- Get reward r
- New state $s' \in \mathcal{S}$



robotics

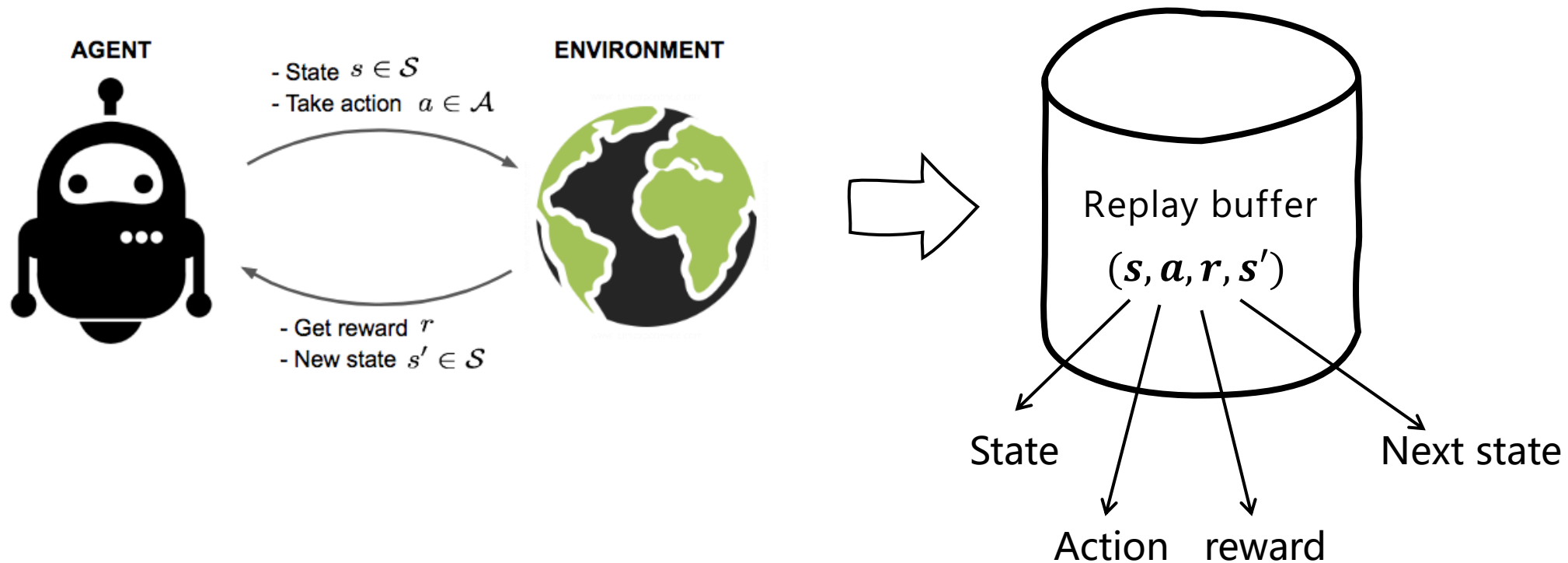


Games

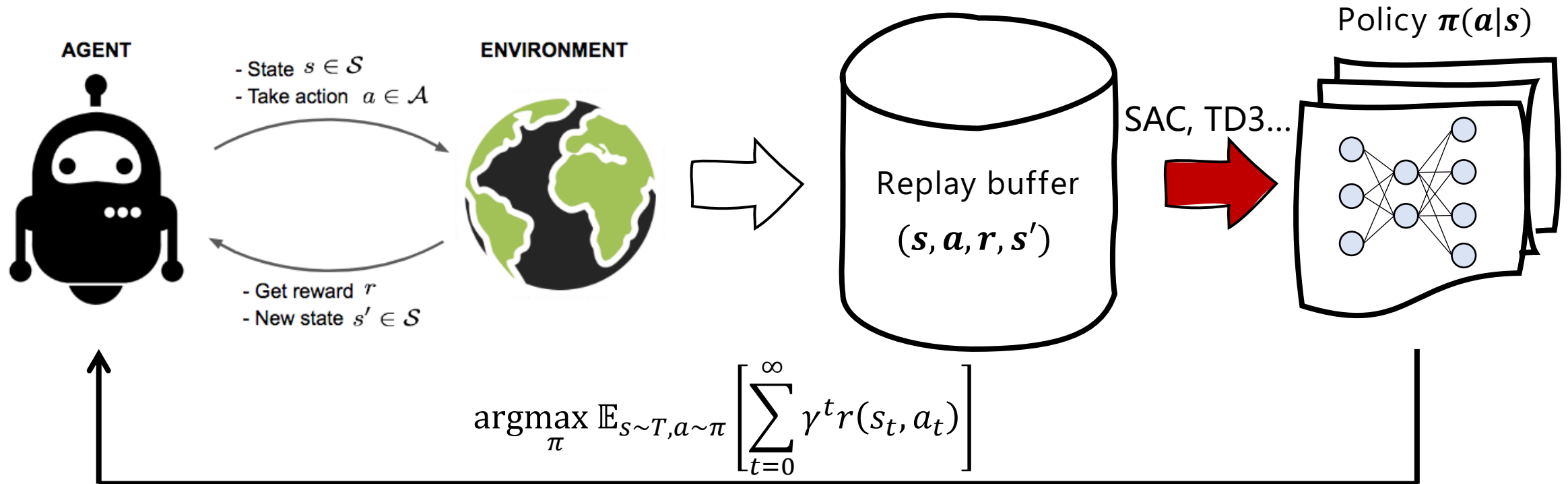


Automatic driving

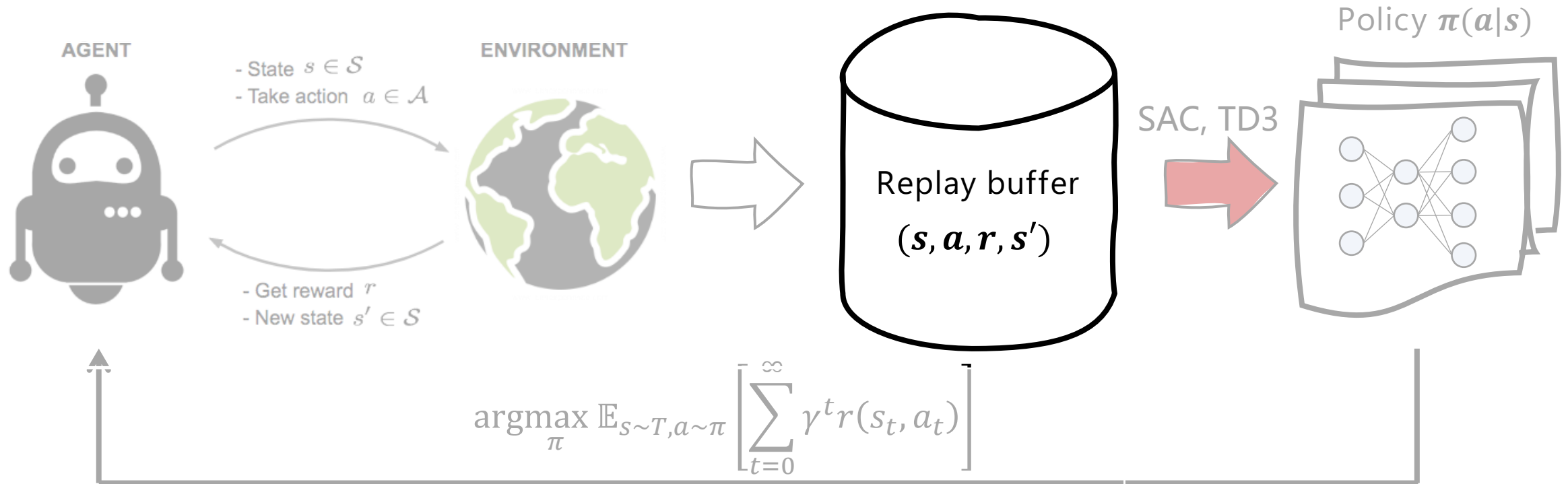
Then collecting data in replay buffer



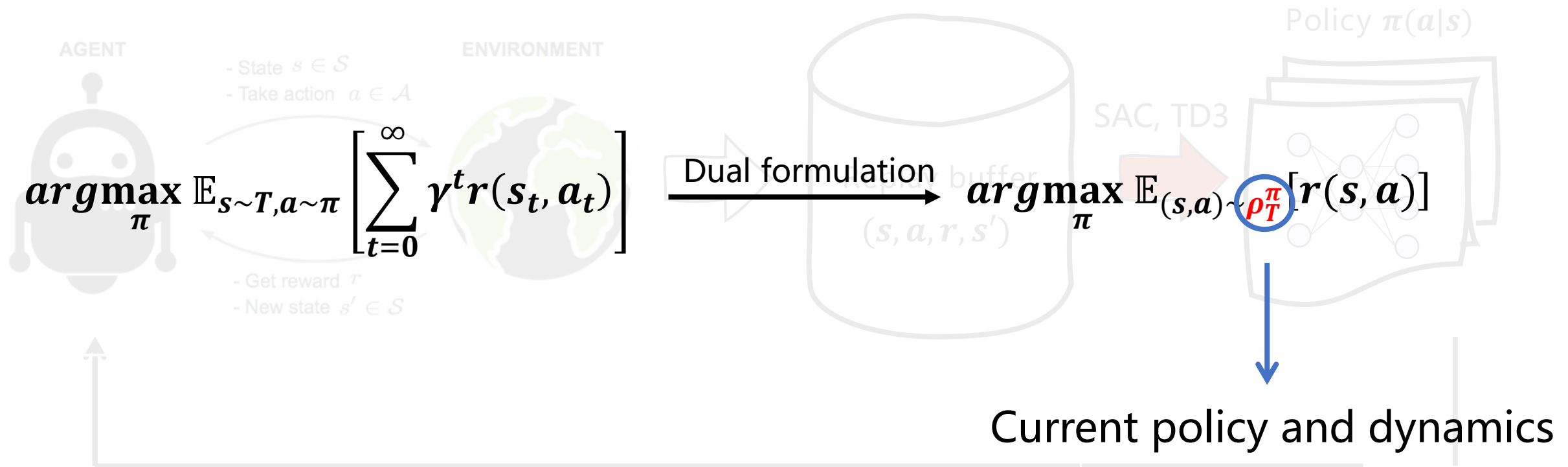
And training policy by RL



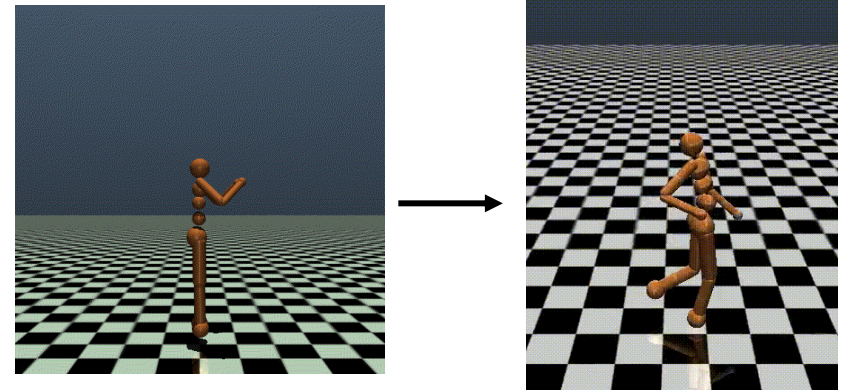
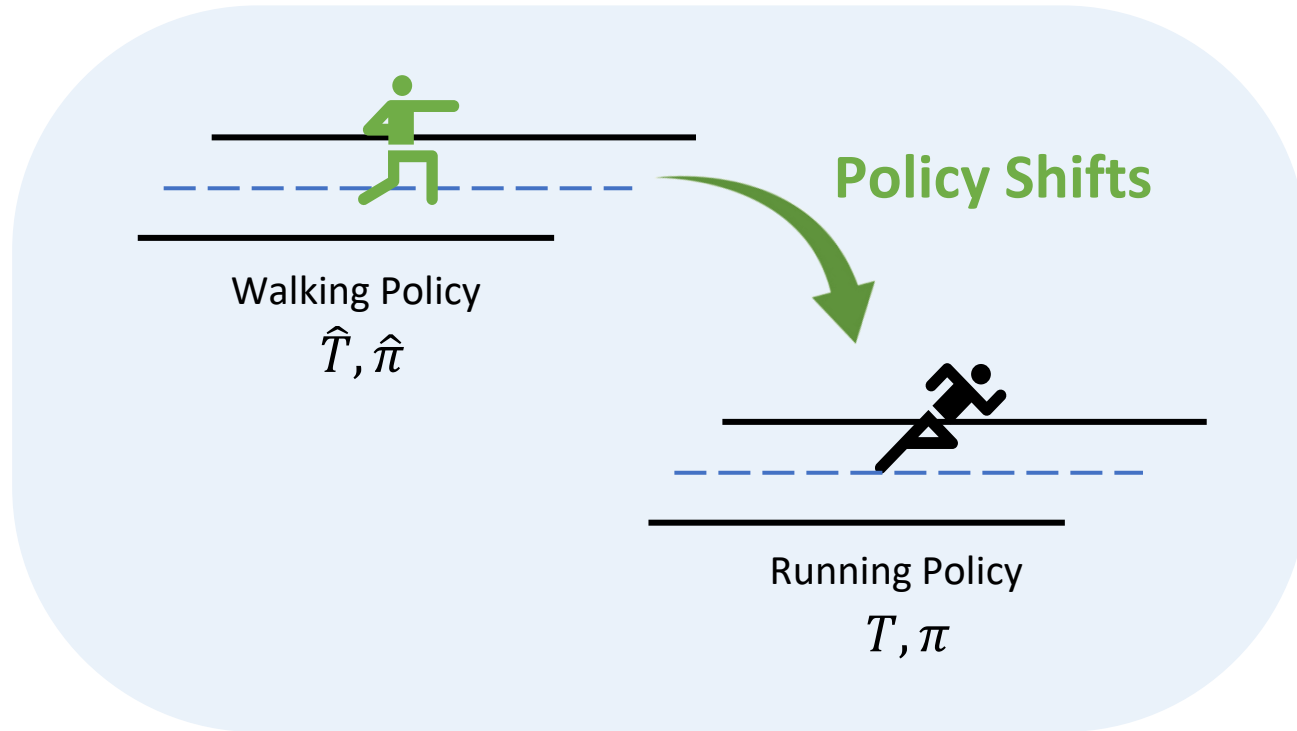
And training policy by RL



And training policy by RL

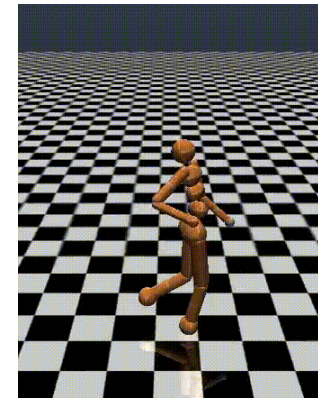
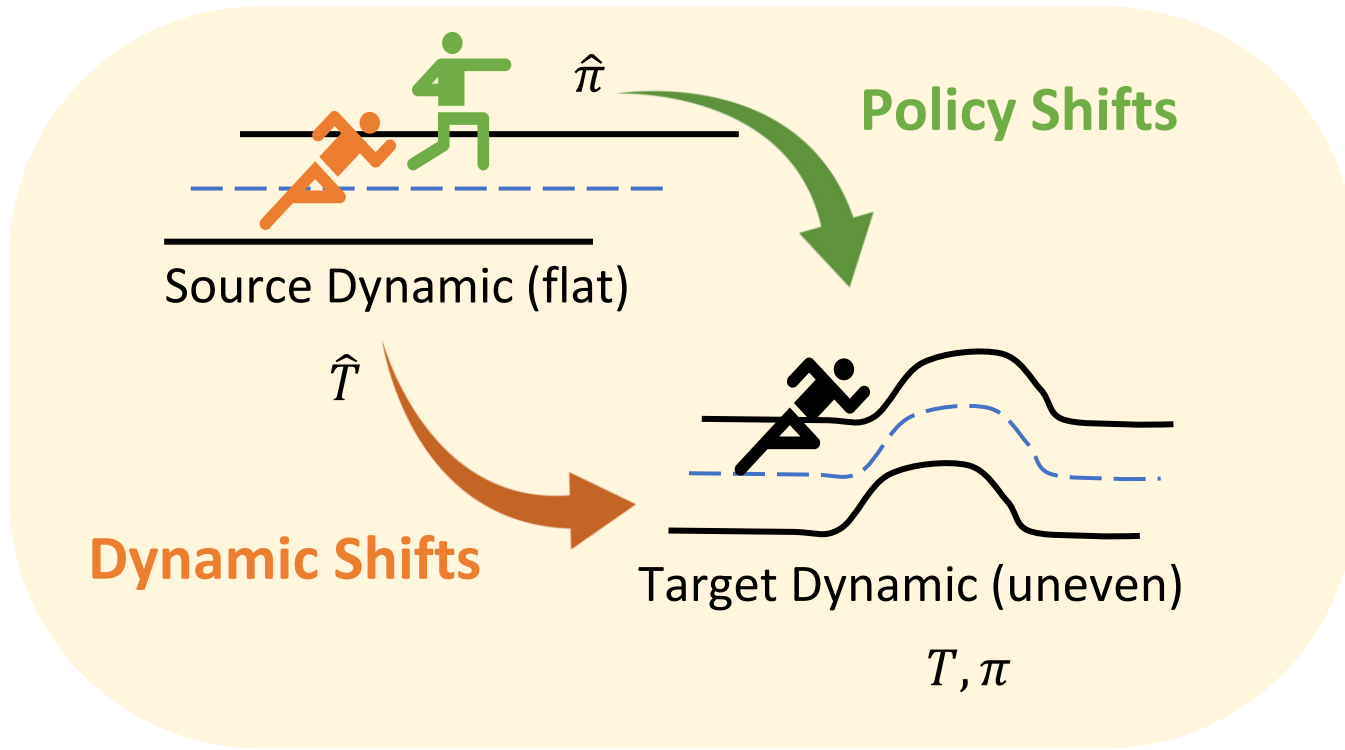


Three cases: Stationary environment

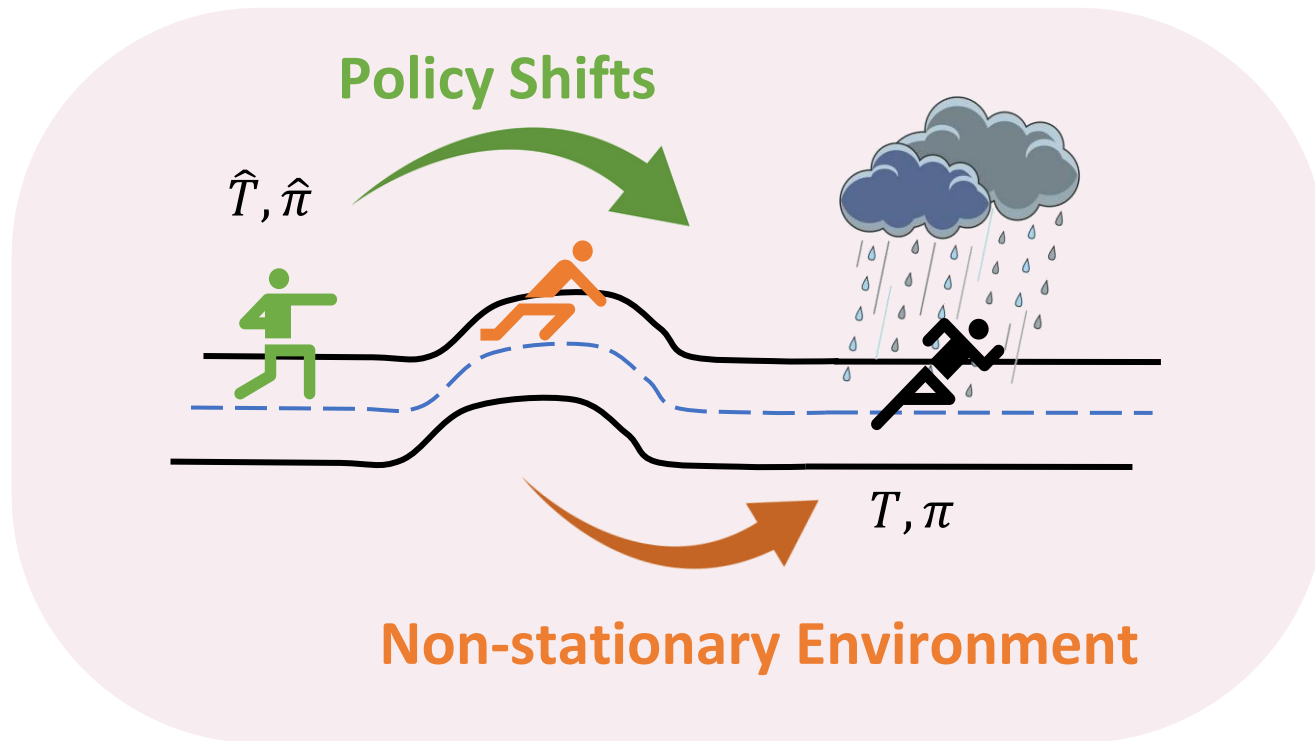


$$T \simeq \hat{T}, \pi \neq \hat{\pi}$$

Three cases: Domain adaption



Three cases: Non-stationary environment



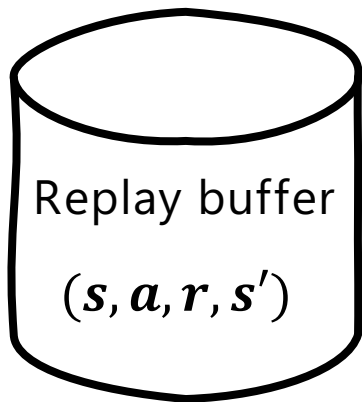
$$T \neq \hat{T}, \pi \neq \hat{\pi}$$



To figure out the distribution gaps

Transition occupancy distribution measure [Viano et al., 2021; Ma et al., 2023]

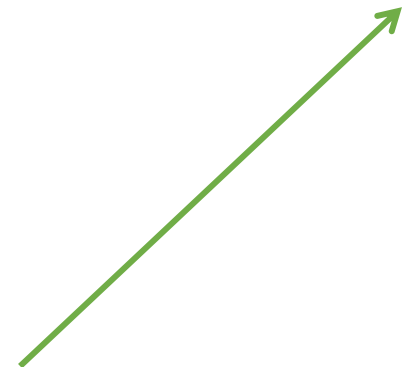
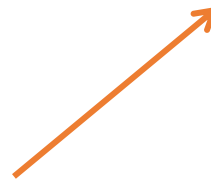
$$\rho_{\hat{T}}^{\hat{\pi}}(s, a, s') = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr[s_t = s, a_t = a, s_{t+1} = s' | s_0 \sim \mu_0, a_t \sim \hat{\pi}(\cdot | s_t), s_{t+1} \sim \hat{T}(\cdot | s_t, a_t)]$$



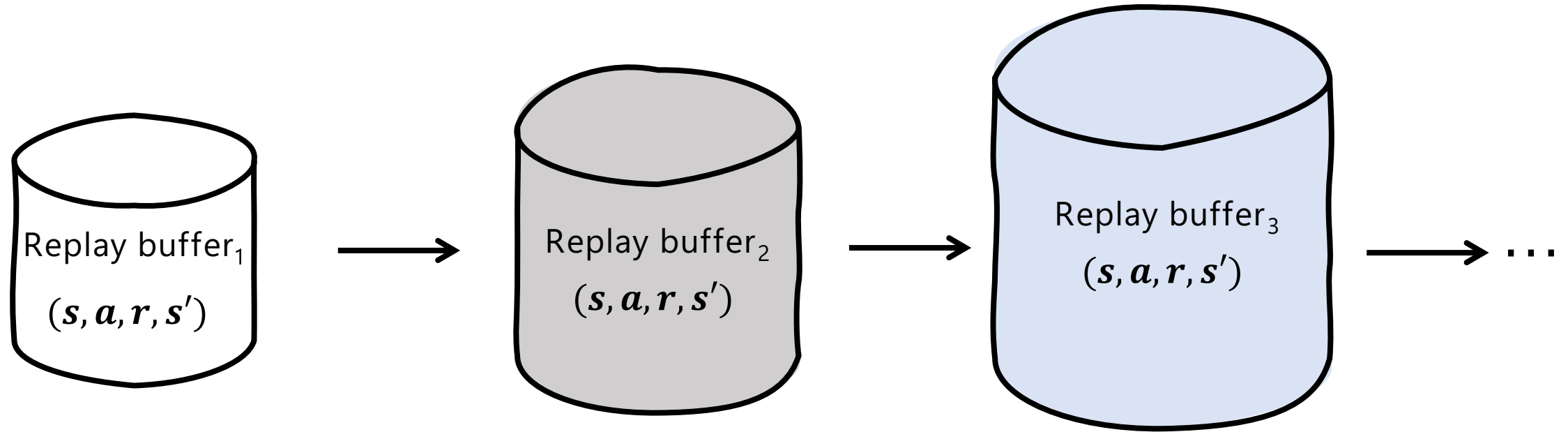
Behavior Policy $\hat{\pi}(a|s)$



Empirical Dynamics $\hat{T}(s'|s, a)$



When data are collected under different π and T



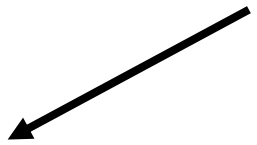
$$\begin{cases} \hat{\pi} \neq \pi \\ \hat{T} \neq T \end{cases} \Rightarrow \rho_{\hat{T}}^{\hat{\pi}}(s, a, s') \neq \rho_T^{\pi}(s, a, s')$$

The influence of distribution gaps in RL

$$\underbrace{J(\pi) = \mathbb{E}_{(s,a) \sim P, \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]}_{\text{General objective}} \quad \Rightarrow \quad \underbrace{J(\pi) = \mathbb{E}_{(s,a,s') \sim \rho_T^\pi [r(s, a)]}_{\text{Dual objective}}}_{\text{Current data distribution}}$$

$$\begin{cases} \hat{\pi} \neq \pi \\ \hat{T} \neq T \end{cases} \quad \Rightarrow \quad \rho_{\hat{T}}^{\hat{\pi}}(s, a, s') \neq \rho_T^\pi(s, a, s')$$

How to learn policy when $\rho_{\hat{T}^{\hat{\pi}}}(s, a, s') \neq \rho_T^{\pi}(s, a, s')$?


$$J(\pi) = \mathbb{E}_{(s,a,s') \sim \rho_T^{\pi}}[r(s, a)]$$

A surrogate objective with distribution shifts

$$J(\pi) = \mathbb{E}_{(s,a,s') \sim \rho_T^\pi} [r(s, a)]$$

A surrogate objective with distribution shifts

$$J(\pi) = \mathbb{E}_{(s,a,s') \sim \rho_T^\pi} [r(s, a)]$$

$$> \log \mathbb{E}_{(s,a,s') \sim \rho_T^\pi} [r(s, a)]$$

A surrogate objective with distribution shifts

$$\begin{aligned} J(\pi) &= \mathbb{E}_{(s,a,s') \sim \rho_T^\pi} [r(s, a)] \\ &> \log \mathbb{E}_{(s,a,s') \sim \rho_T^\pi} [r(s, a)] \\ &= \log \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^\pi} \left[\left(\frac{\rho_T^\pi}{\rho_{\hat{T}}^\pi} \right) r(s, a) \right] \end{aligned}$$

A surrogate objective with distribution shifts

$$\begin{aligned} J(\pi) &= \mathbb{E}_{(s,a,s') \sim \rho_T^\pi} [r(s, a)] \\ &> \log \mathbb{E}_{(s,a,s') \sim \rho_T^\pi} [r(s, a)] \\ &= \log \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^\pi} \left[\left(\frac{\rho_T^\pi}{\rho_{\hat{T}}^\pi} \right) r(s, a) \right] \\ &\geq \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^\pi} \left[\log \left(\frac{\rho_T^\pi}{\rho_{\hat{T}}^\pi} \right) + \log r(s, a) \right] \end{aligned}$$

A surrogate objective with distribution shifts

$$\begin{aligned} J(\pi) &= \mathbb{E}_{(s,a,s') \sim \rho_T^\pi} [r(s, a)] \\ &> \log \mathbb{E}_{(s,a,s') \sim \rho_T^\pi} [r(s, a)] \\ &= \log \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^\pi} \left[\left(\frac{\rho_T^\pi}{\rho_{\hat{T}}^\pi} \right) r(s, a) \right] \\ &\geq \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^\pi} \left[\log \left(\frac{\rho_T^\pi}{\rho_{\hat{T}}^\pi} \right) + \log r(s, a) \right] \\ &= \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^\pi} [\log r(s, a)] - D_{\text{KL}}(\rho_{\hat{T}}^\pi \parallel \rho_T^\pi) \end{aligned}$$

A surrogate objective with distribution shifts

$$\hat{J}(\pi) = \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^{\pi}} \left[\log r(s, a) - \alpha \log \left(\frac{\rho_T^{\pi}}{\rho_{\hat{T}}^{\pi}} \right) \right] - \alpha D_f(\rho_{\hat{T}}^{\pi} \| \rho_T^{\pi})$$

policy & dynamics shifts

policy shifts

$$s.t. \quad \rho^{\pi}(s, a) = (1 - \gamma)\mu_0(s)\pi(a|s) + \gamma \mathcal{T}_*^{\pi} \rho^{\pi}(s, a)$$

A surrogate objective with distribution shifts

$$\hat{J}(\pi) = \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^{\pi}} \left[\log r(s, a) - \alpha \log \left(\frac{\rho_T^{\pi}}{\rho_{\hat{T}}^{\pi}} \right) \right] - \alpha D_{\text{KL}}(\rho_{\hat{T}}^{\pi} \| \rho_T^{\pi})$$

policy & dynamics shifts

policy shifts

$$s.t. \quad \rho^{\pi}(s, a) = (1 - \gamma)\mu_0(s)\pi(a|s) + \gamma \mathcal{T}_*^{\pi} \rho^{\pi}(s, a)$$

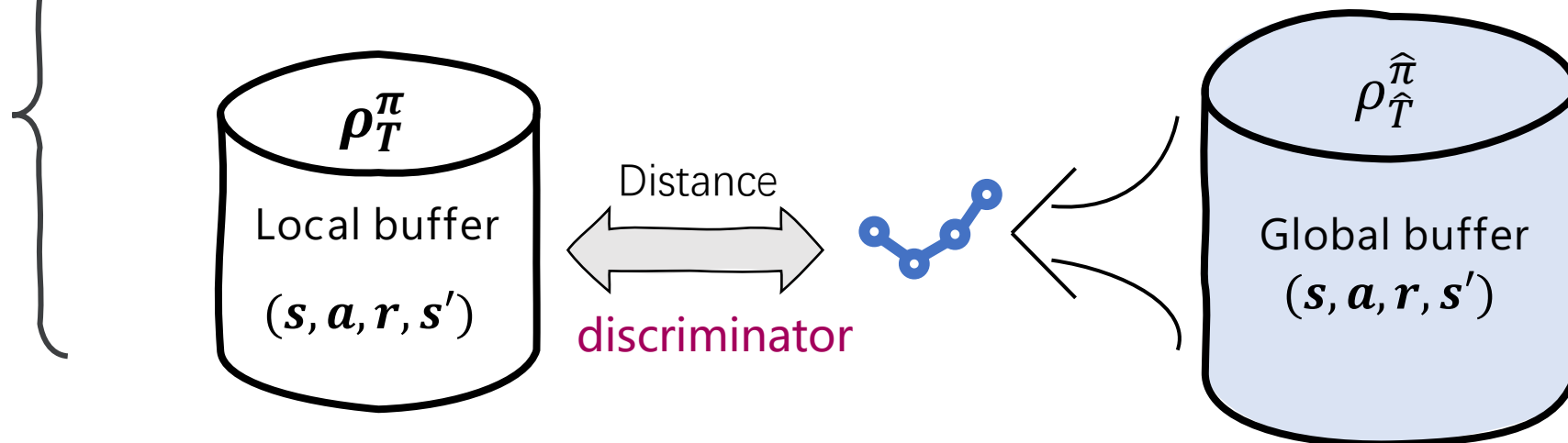
How to solve the constrained optimization problem?

Three theoretical steps for tractable solution

$$\hat{J}(\pi) = \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^{\pi}} \left[\log r(s, a) - \alpha \log \left(\frac{\rho_T^{\pi}}{\rho_{\hat{T}}^{\pi}} \right) \right] - \alpha D_f(\rho_{\hat{T}}^{\pi} \| \rho_T^{\pi})$$

$$s.t. \quad \rho^{\pi}(s, a) = (1 - \gamma)\mu_0(s)\pi(a|s) + \gamma \mathcal{T}_*^{\pi} \rho^{\pi}(s, a)$$

Step1: Deal with the discrepancy $\log \left(\frac{\rho_T^{\pi}}{\rho_{\hat{T}}^{\pi}} \right)$ by a discriminator



Three theoretical steps for tractable solution

$$\hat{J}(\pi) = \mathbb{E}_{(s,a,s') \sim \rho_{\hat{T}}^{\pi}} \left[\log r(s, a) - \alpha \log \left(\frac{\rho_T^{\pi}}{\rho_{\hat{T}}^{\pi}} \right) \right] - \alpha D_f(\rho_{\hat{T}}^{\pi} \| \rho_T^{\pi})$$

$$s.t. \quad \rho^{\pi}(s, a) = (1 - \gamma)\mu_0(s)\pi(a|s) + \gamma \mathcal{T}_*^{\pi} \rho^{\pi}(s, a)$$

Step1: Deal with the discrepancy $\log \left(\frac{\rho_T^{\pi}}{\rho_{\hat{T}}^{\pi}} \right)$ by a discriminator

Step2: Deal with the constraint $\rho^{\pi}(s, a)$ by Lagrange multipliers

Step3: Deal with the unknown term $\rho_{\hat{T}}^{\pi}(s, a, s')$ by Fenchel conjugate

Our solution

$$\min_{\pi} \max_{Q(s,a)} \underbrace{(1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi} [Q(s, a)]}_{\text{Initial distribution}} + \alpha \mathbb{E}_{(s,a,s') \sim \rho_{\hat{\pi}}^T} \left[f_{\star} \left(\frac{\log r(s, a) - R(s, a, s') + \gamma \mathcal{T}^{\pi} Q(s, a) - Q(s, a)}{\alpha} \right) \right]$$

Fenchel conjugate function

Sample from collected data

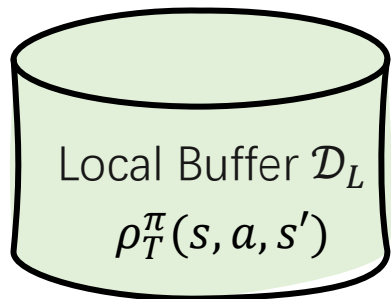
Distribution gap between
Current data &
Historical data

TD error

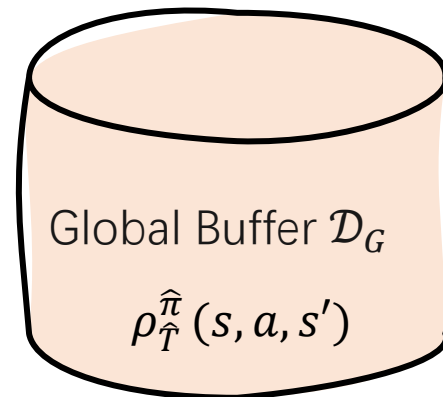
Algorithm: OMPO

$$\min_{\pi} \max_{Q(s,a)} (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi} [Q(s, a)] + \alpha \mathbb{E}_{(s,a,s') \sim \rho_{\hat{\pi}}^{\pi}} \left[f_{\star} \left(\frac{\log r(s, a) - R(s, a, s') + \gamma \mathcal{T}^{\pi} Q(s, a) - Q(s, a)}{\alpha} \right) \right]$$

- Critic network for $Q(s, a)$, Inner loop optimization
- Policy network for $\pi(a|s)$, Outer loop optimization



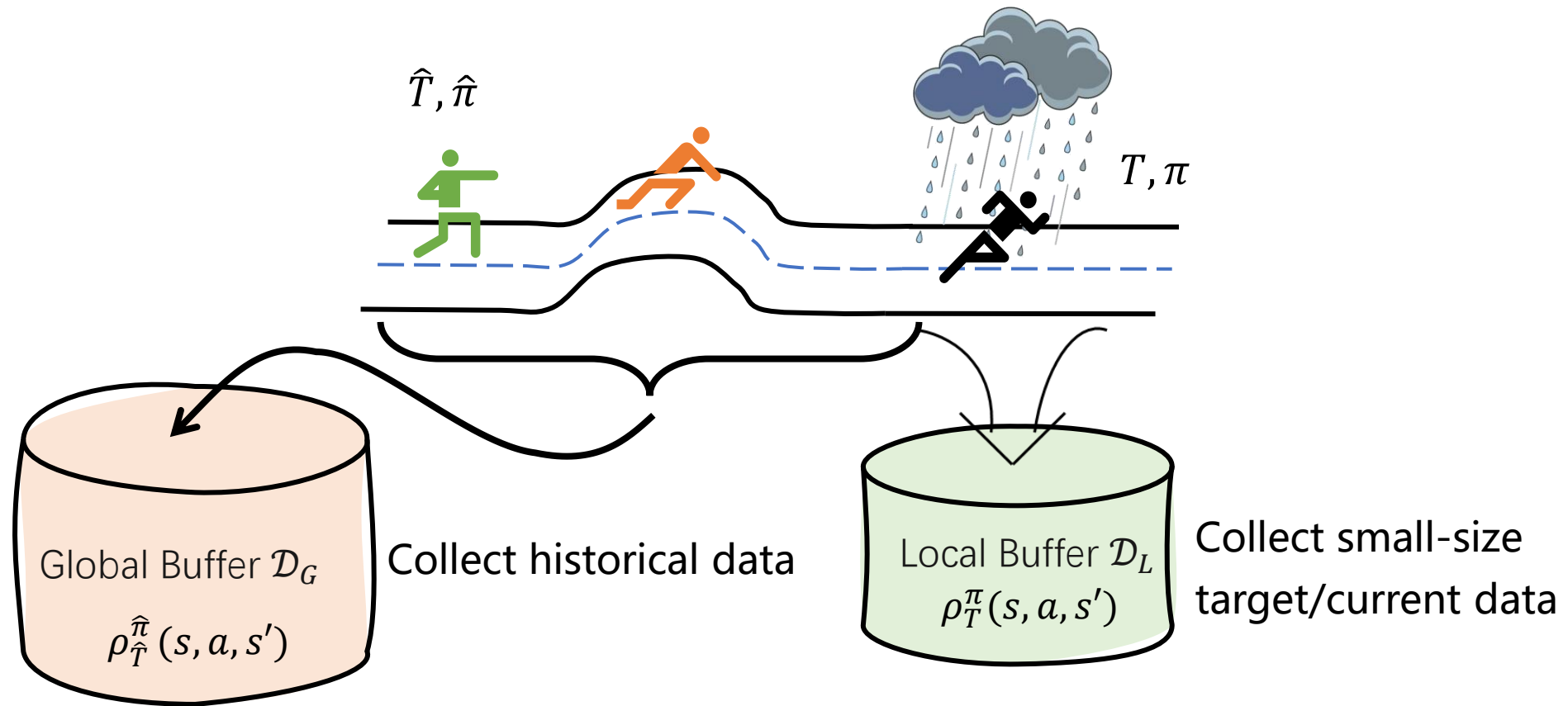
Collect target/current data



Collect historical data

An example of how to the two buffer

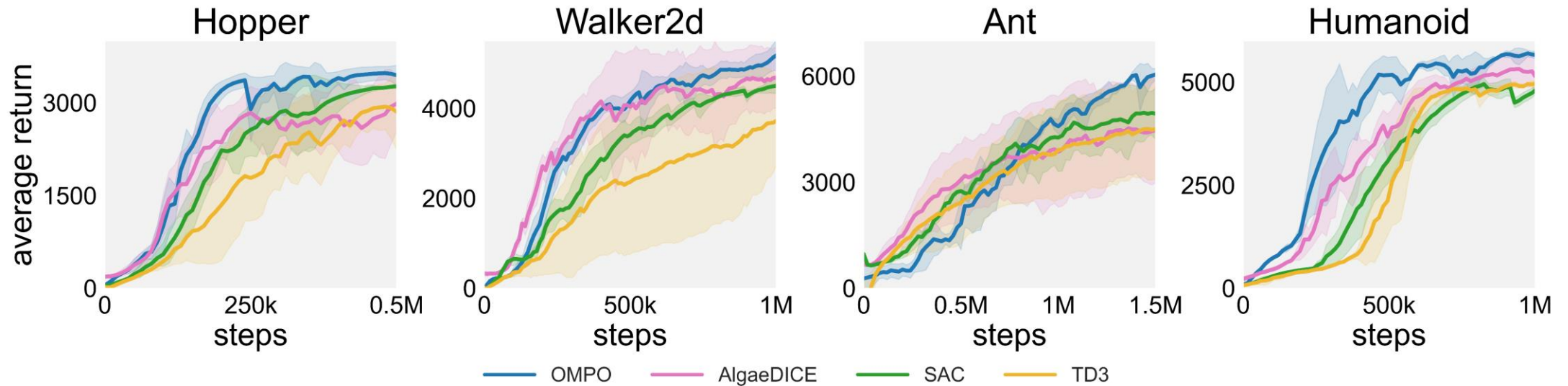
- Non-stationary environments



Experiments

Three scenarios with specialised baselines

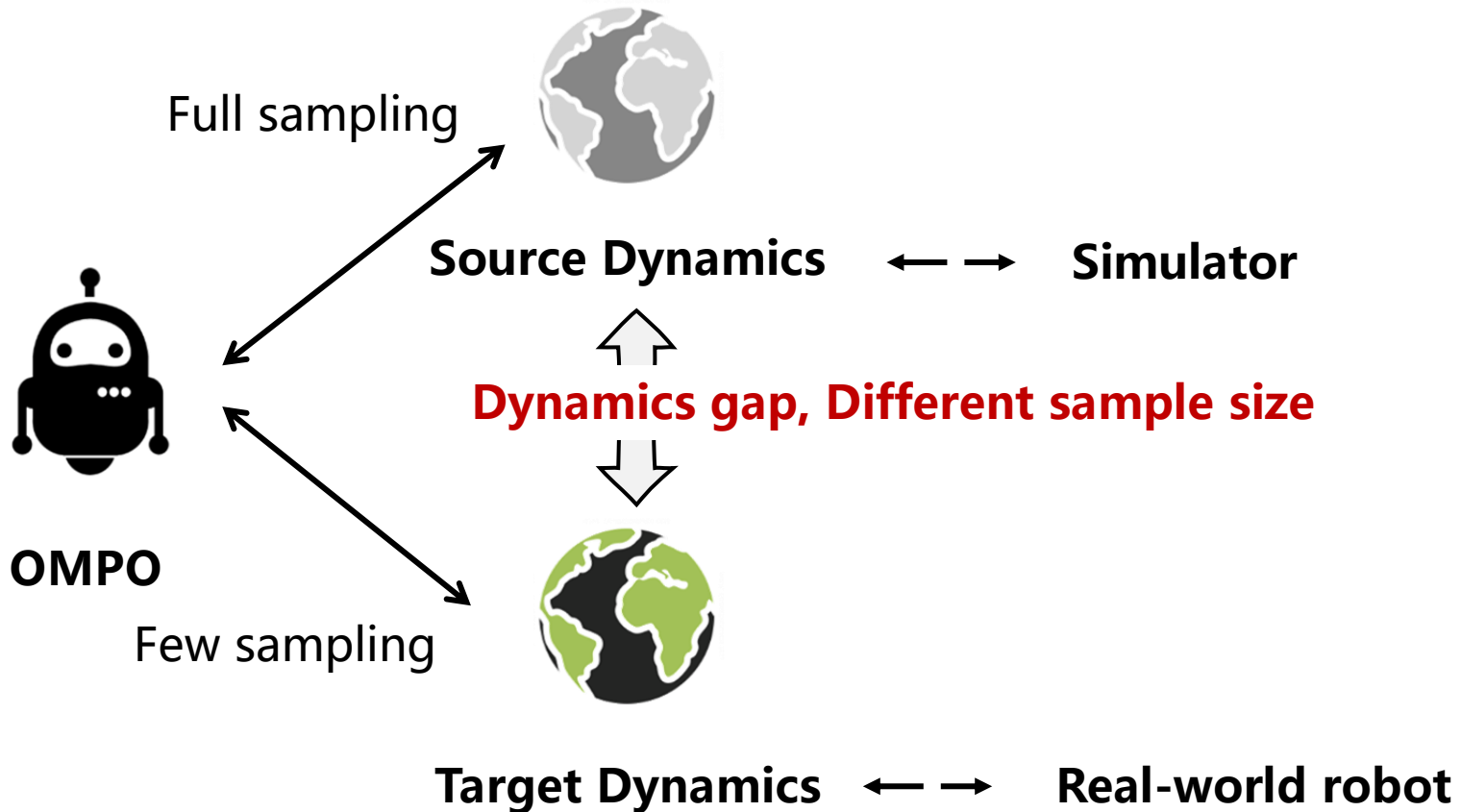
Stationary environment: Hopper, Walker2d, Ant, Humanoid



Stable performance & sample efficiency

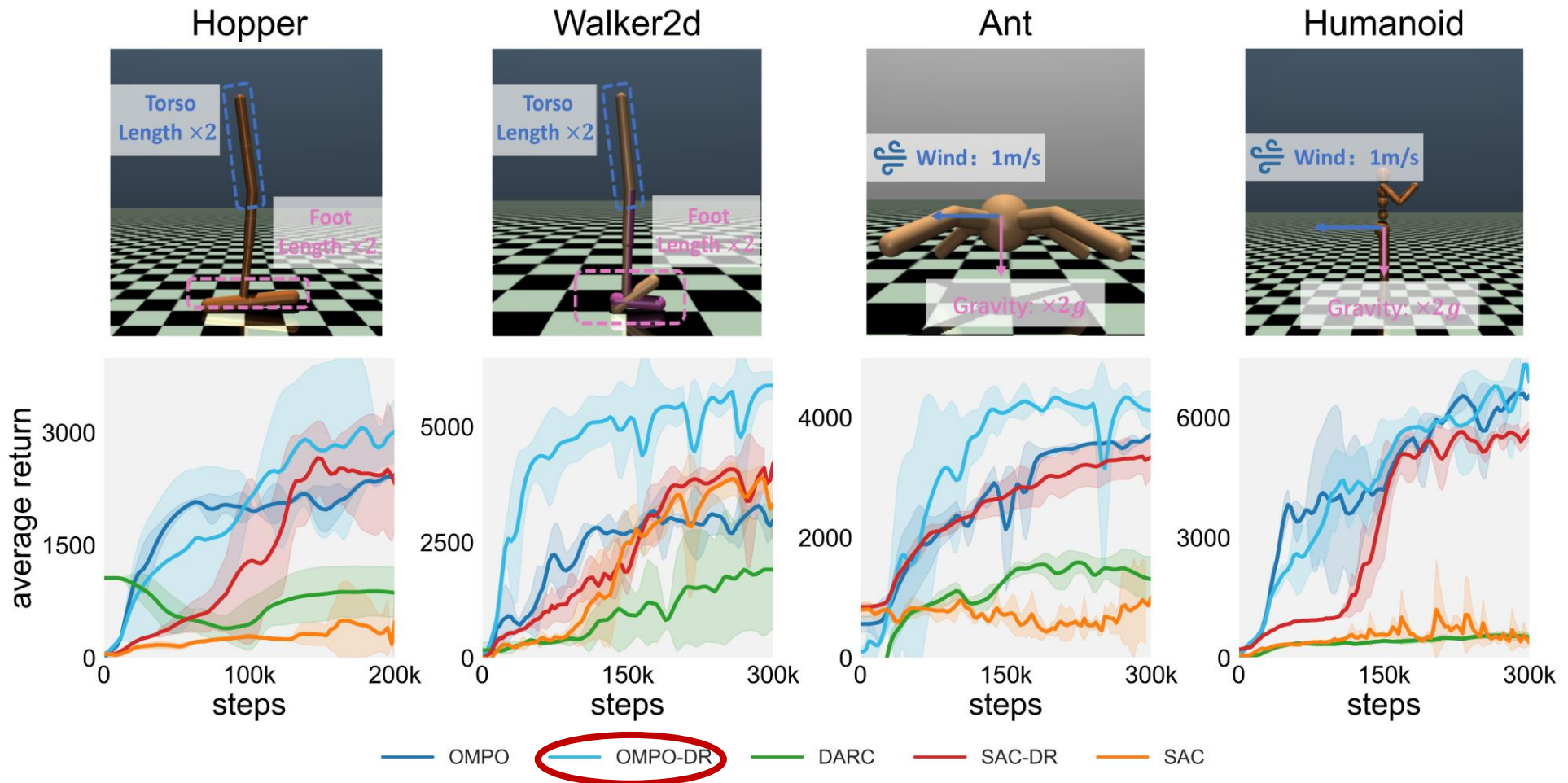
Three scenarios with specialised baselines

Domain Adaption: Hopper, Walker2d, Ant, Humanoid



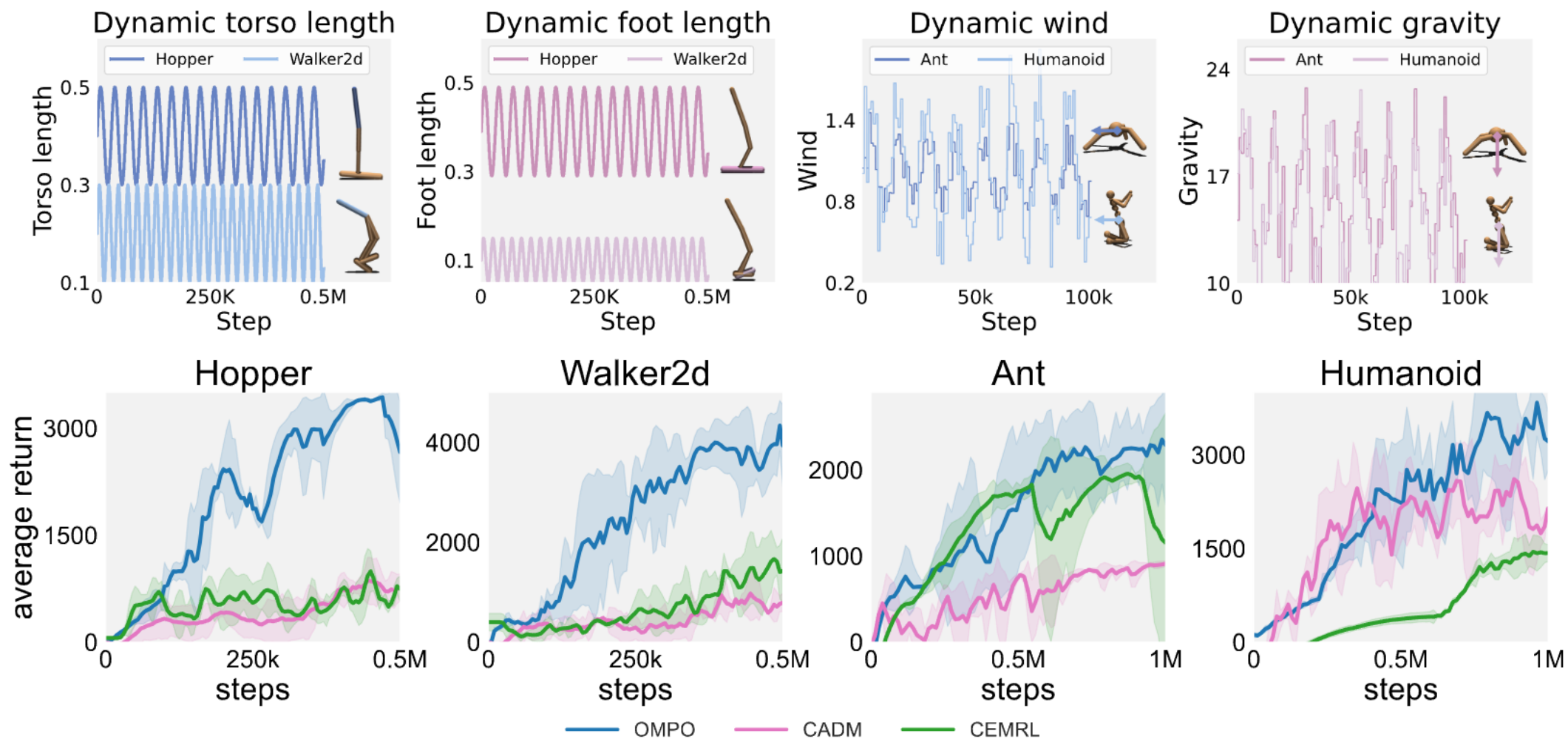
Three scenarios with specialised baselines

Domain Adaption: Hopper, Walker2d, Ant, Humanoid



Three scenarios with specialised baselines

Non-stationary environment: Hopper, Walker2d, Ant, Humanoid



Conclusion

- We propose a general surrogate objective for policy and dynamics shifts
- We develop a unified framework to tackle diverse shift settings

Future research

- The two-buffer setting can be extended, like offline-to-online RL, hybrid RL and imitation learning
- How to adaptively decide the local buffer size would be interesting

Thank you for listening!

