



POLITECNICO  
MILANO 1863

# POLICY OPTIMIZATION AS ONLINE LEARNING WITH MEDIATOR FEEDBACK

ALBERTO MARIA METELLI\*, MATTEO PAPINI\*, PIERLUCA D'ORO, AND MARCELLO RESTELLI

{albertomaria.metelli, matteo.papini, marcello.restelli}@polimi.it, pierluca.doro@mail.polimi.it



## MOTIVATION AND IDEA

**Problem:** How to deal with **exploration** in **Policy Optimization (PO)**?

**Idea:** exploit the **inherent structure** of the PO problem via **multiple importance sampling**

## POLICY OPTIMIZATION

- **Parameter space**  $\Theta \subseteq \mathbb{R}^d$
- A parametric **policy** for each  $\theta \in \Theta$
- Each inducing a distribution  $p_\theta$  over **trajectories**
- A **return**  $\mathcal{R}(\tau)$  for every trajectory  $\tau$
- **Goal:** maximize the **expected return** (Deisenroth et al., 2013)

$$\theta^* \in \arg \max_{\theta \in \Theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [\mathcal{R}(\tau)]$$

## POLICY OPTIMIZATION AS ONLINE LEARNING

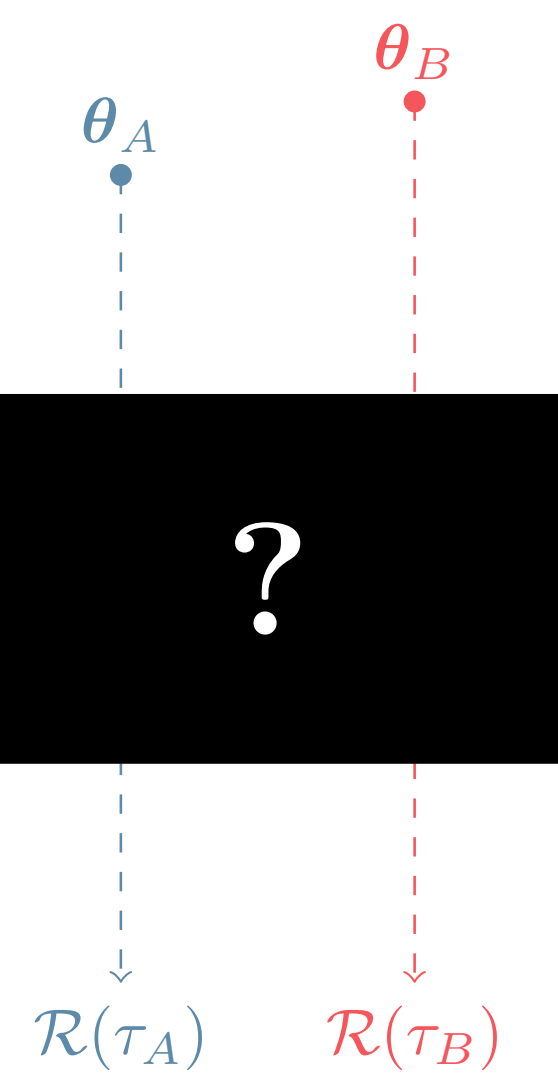
For  $t = 1, 2, \dots$

- **Select** parameter  $\theta_t$  and run  $\pi_{\theta_t}$
- **Observe** the trajectory  $\tau_t$  and the return  $\mathcal{R}(\tau_t)$

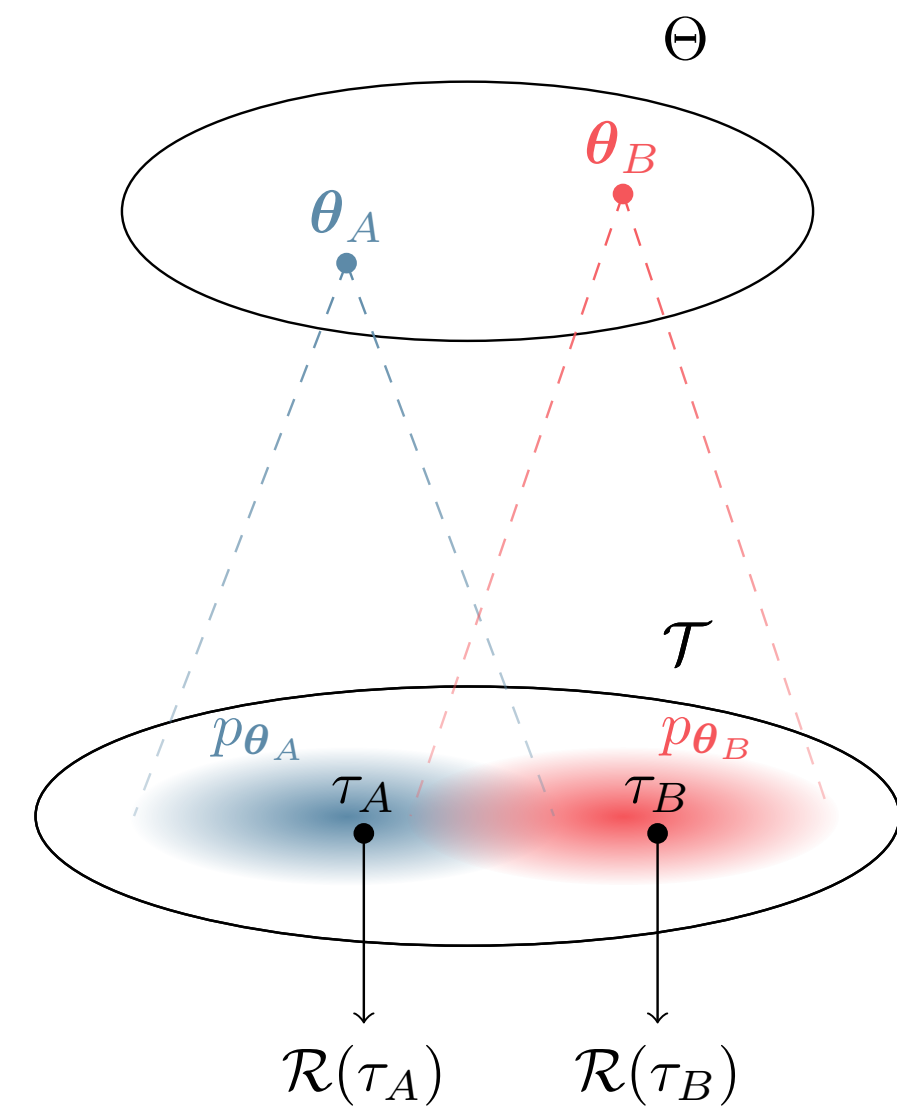
**Goal:** minimize the **regret** (Auer et al., 2002)

$$\text{Regret}(n) = \sum_{t=1}^n J(\theta^*) - J(\theta_t) = \sum_{t=1}^n \Delta(\theta_t)$$

**Bandit Feedback**



**Mediator Feedback**



## REGRET LOWER BOUNDS

- Two-parameter space  $\Theta = \{\theta_A, \theta_B\}$
- Performance gap  $\Delta = J(\theta_A) - J(\theta_B)$

- If  $D_{KL}(p_{\theta_A} \| p_{\theta_B}) < \infty$  and  $D_{KL}(p_{\theta_B} \| p_{\theta_A}) < \infty \implies$  **constant regret**

$$\mathbb{E} \text{Regret}(n) \geq \mathcal{O}\left(\frac{1}{\Delta}\right)$$

- If  $D_{KL}(p_{\theta_A} \| p_{\theta_B}) = \infty$  or  $D_{KL}(p_{\theta_B} \| p_{\theta_A}) = \infty \implies$  **logarithmic regret**

$$\mathbb{E} \text{Regret}(n) \geq \mathcal{O}\left(\frac{1}{\Delta} \log(\Delta^2 n)\right)$$

## IMPORTANCE SAMPLING FOR MEDIATOR FEEDBACK

- **Idea:** use **all** the samples to estimate the expected return of **any** policy  $\mathcal{H}_t = \{(\theta_i, \tau_i, \mathcal{R}(\tau_i))\}_{i \in [t-1]}$
- **Mixture** distribution  $\Phi_t$  of the **behavioral policies** played and relative **multiple importance weight** with *balance heuristic* (Veach and Guibas, 1995) w.r.t. the **target policy**  $p_\theta(\tau_i)$ :

$$\Phi_t = \frac{1}{t-1} \sum_{j=1}^{t-1} p_{\theta_j}(\tau_i) \implies \frac{p_\theta(\tau_i)}{\Phi_t(\tau_i)} = \frac{p_\theta(\tau_i)}{\frac{1}{t-1} \sum_{j=1}^{t-1} p_{\theta_j}(\tau_i)} = \frac{\prod_{h=0}^{H-1} \pi_\theta(a_{ih} | s_{ih})}{\frac{1}{t-1} \sum_{j=1}^{t-1} \prod_{h=0}^{H-1} \pi_{\theta_j}(a_{ih} | s_{ih})}$$

- Vanilla importance weight leads to **heavy-tailed** estimator (Metelli et al., 2018)  $\implies$  employ a *time-variant* weight **truncation** threshold  $M_t(\theta)$  (Ionides, 2008)

$$\check{J}_t(\theta) = \frac{1}{t-1} \sum_{i=1}^{t-1} \min \left\{ \frac{p_\theta(\tau_i)}{\Phi_t(\tau_i)}, M_t(\theta) \right\} \mathcal{R}(\tau_i)$$

$$M_t(\theta) = \sqrt{\frac{(t-1) d_2(p_\theta \| \Phi_t)}{\log \frac{1}{\delta}}}$$

**Truncation threshold**

$$d_2(p_\theta \| \Phi_t) = \int \frac{p_\theta(\tau_i)^2}{\Phi_t(\tau_i)} d\tau$$

**Renyi divergence**

$$\check{J}_t(\theta) - J(\theta) \leq 2.75 \sqrt{\frac{\log \frac{1}{\delta}}{\eta_t(\theta)}}$$

$$\eta_t(\theta) = \frac{t-1}{d_2(p_\theta \| \Phi_t)}$$

**Effective sample size**

- We obtain **exponential concentration** (Papini et al., 2019; Metelli et al., 2020):

## ALGORITHMS

Execute  $\pi_{\theta_1}$ , observe  $\tau_1 \sim p_{\theta_1}$  and  $\mathcal{R}(\tau_1)$

**for**  $t = 2, \dots, n$  **do**

    Compute expected return estimate  $\check{J}_t(\theta)$

    Select  $\theta_t \in \arg \max_{\theta \in \Theta} B_t(\theta)$

    Execute  $\pi_{\theta_t}$ , observe  $\tau_t \sim p_{\theta_t}$  and  $\mathcal{R}(\tau_t)$

**end for**

**OPTIMIST** (Papini et al., 2019)

Compute upper confidence bound:

$$B_t(\theta) = \check{J}_t(\theta) + 2.42 \sqrt{\frac{\alpha \log t}{\eta_t(\theta)}}$$

**RANDOMIST** (new!)

Generate perturbation:

$$U_t(\theta) = \frac{1}{\eta_t(\theta)} \sum_{l=1}^{a\eta_t(\theta)} \tau_l + b, \text{ with } \tau_l \sim \text{Ber}(1/2)$$

Compute index  $B_t(\theta) = \check{J}_t(\theta) + U_t(\theta)$

## REGRET UPPER BOUNDS COMPARISON

Finite Policy Space

$$v = \max_{\theta, \theta' \in \Theta} d_2(p_\theta \| p_{\theta'}) \text{ and } \Delta = \min_{\theta \neq \theta^*} J(\theta^*) - J(\theta)$$

Algorithm	Exploration	$\mathbb{E} \text{Regret}(n)$	
		$v = \infty$	$v < \infty$
Greedy	-	$\mathcal{O}(n)$	$\mathcal{O}\left(\frac{v}{\Delta} \log \frac{v}{\Delta^2}\right)$
UCB1	deterministic	$\mathcal{O}\left(\frac{1}{\Delta} \log n\right)$	$\mathcal{O}\left(\frac{1}{\Delta} \log n\right)$
OPTIMIST	deterministic	$\mathcal{O}\left(\frac{1}{\Delta} \log n\right)$	$\mathcal{O}\left(\frac{v}{\Delta} \log \frac{v}{\Delta^2}\right)$
RANDOMIST	randomized	$\mathcal{O}\left(\frac{1}{\Delta} \log n\right)$	$\mathcal{O}\left(\frac{v}{\Delta} \log \frac{v}{\Delta^2}\right)$

Compact Policy Space

$$\Theta = [-D, D]^d \text{ and } v = \sup_{\theta, \theta' \in \Theta} d_2(p_\theta \| p_{\theta'})$$

Algorithm	Complexity	$\mathbb{E} \text{Regret}(n)$
OPTIMIST	$t^{1+\frac{d}{2}}$	$\mathcal{O}\left(\sqrt{v d n}\right)$
RANDOMIST	$d t^2$	?

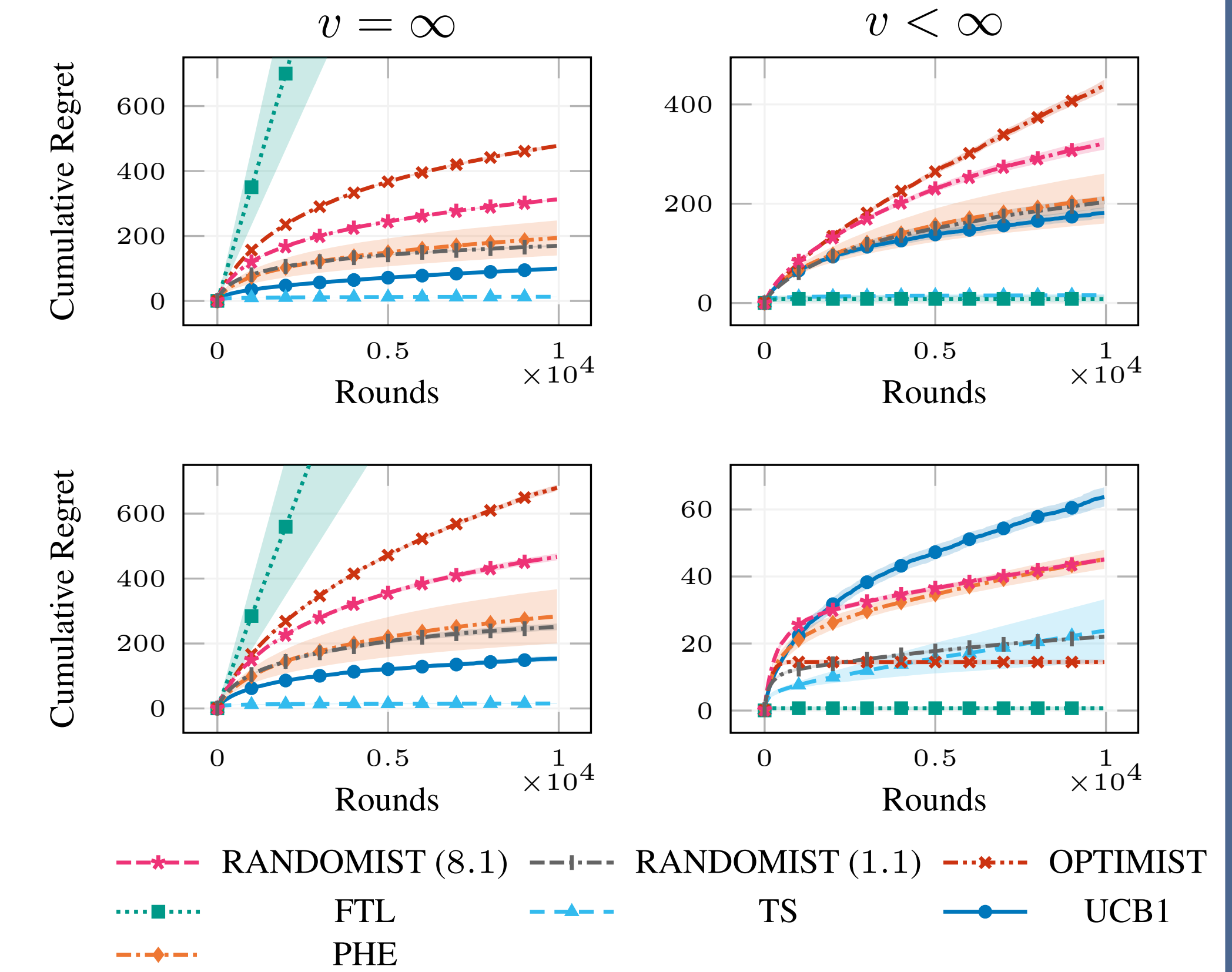
RANDOMIST replaces the **discretization** of OPTIMIST with **MCMC sampling**:

$$\theta_t \sim \Pr \left( \check{J}_t(\theta) + U_t(\theta) = \sup_{\theta' \in \Theta} \check{J}_t(\theta') + U_t(\theta') \right)$$

## EXPERIMENTS

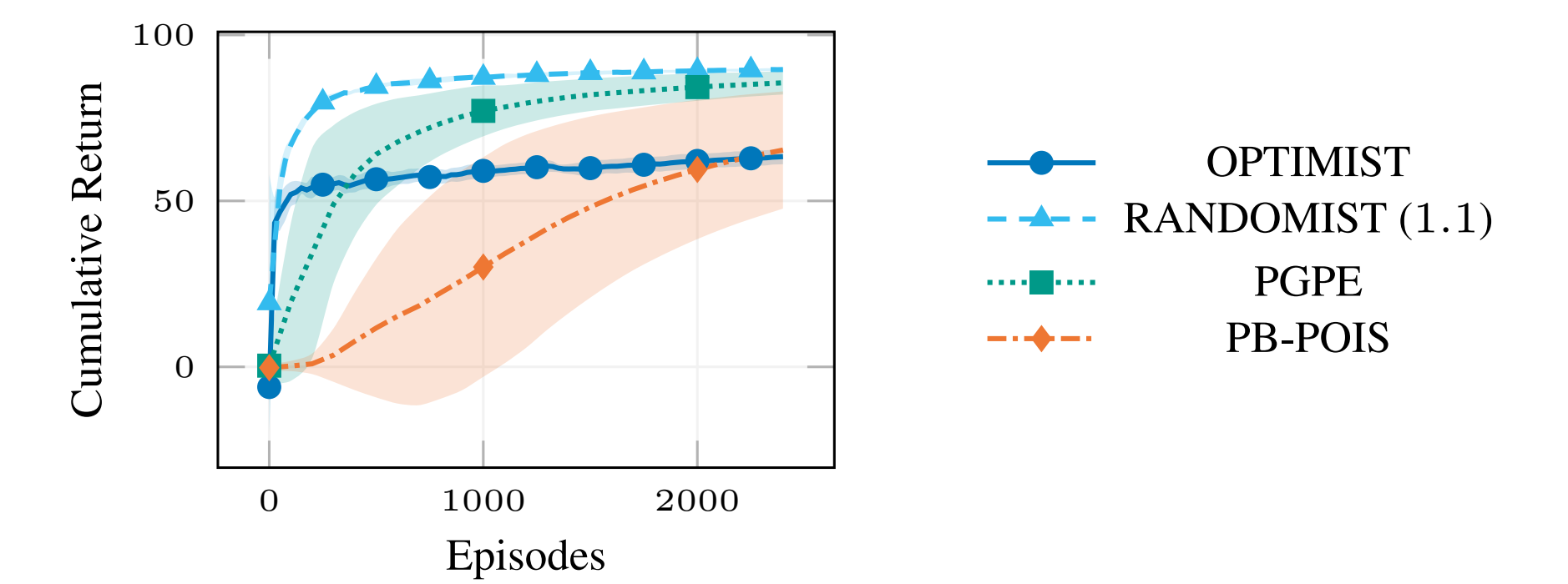
Finite Policy Spaces

**Illustrative Examples (Regret)**

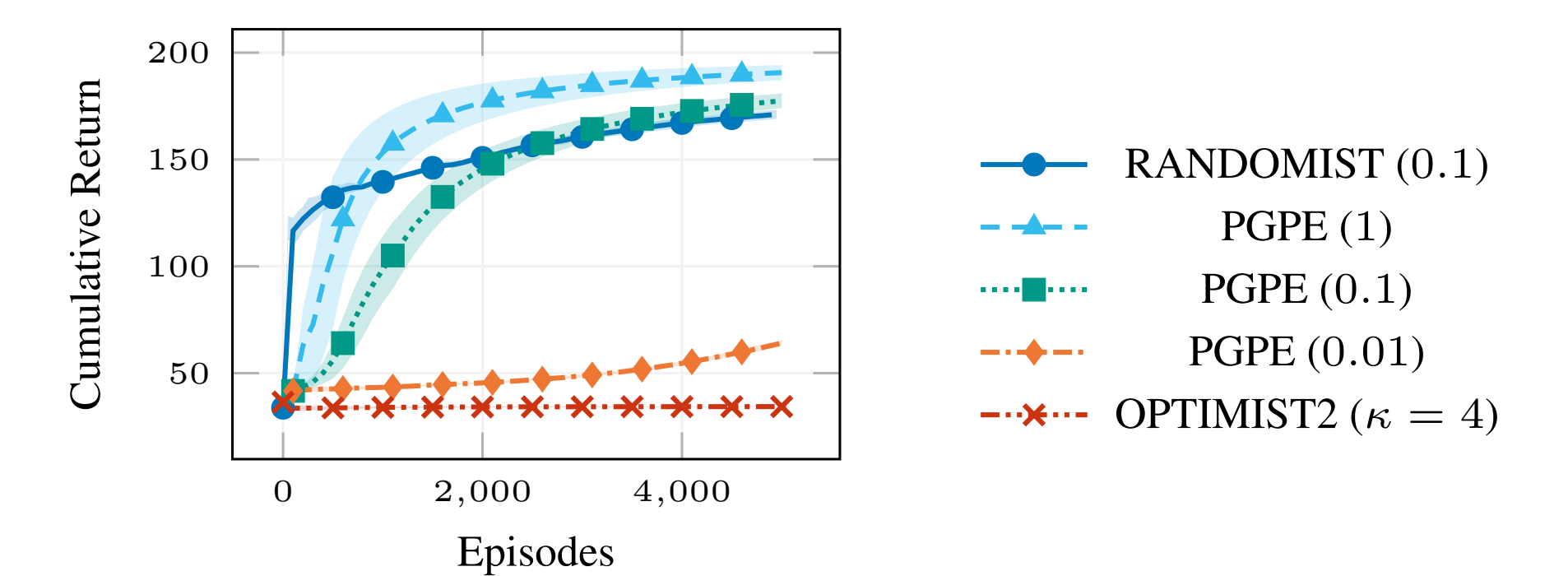


Compact Policy Spaces

**Mountain Car**



**Cartpole**



## REFERENCES

- P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multi-armed bandit problem. *Mach. Learn.*, 47(2-3):235–256, 2002.
- M. P. Deisenroth, G. Neumann, and J. Peters. A survey on policy search for robotics. *Foundations and Trends in Robotics*, 2(1-2):1–142, 2013.
- E. L. Ionides. Truncated importance sampling. *JCGS*, 17(2):295–311, 2008.
- A. M. Metelli, M. Papini, F. Faccio, and M. Restelli. Policy optimization via importance sampling. In *NeurIPS*, 2018.
- A. M. Metelli, M. Papini, N. Montali, and M. Restelli. Importance sampling techniques for policy optimization. *JMLR*, 21(141):1–75, 2020.
- M. Papini, A. M. Metelli, L. Lupo, and M. Restelli. Optimistic policy optimization via multiple importance sampling. In *ICML*, 2019.
- E. Veach and L. J. Guibas. Optimally combining sampling techniques for monte carlo rendering. In S. G. Mair and R. Cook, editors, *SIGGRAPH*, pages 419–428. ACM, 1995.