



POLITECNICO
MILANO 1863

SAFE EXPLORATION IN GAUSSIAN POLICY GRADIENT

MATTEO PAPINI, ANDREA BATTISTELLO AND MARCELLO RESTELLI

matteo.papini@polimi.it, andrea.battistello1994@gmail.com, marcello.restelli@polimi.it



PROBLEM

- **Reinforcement Learning** for **continuous** control [Deisenroth et al., 2013]
- **Policy Gradient (PG)**: iteratively update **parametric** policy π_θ via **gradient ascent** on performance $J(\theta)$ (expected cumulative reward):

$$\theta_{t+1} \leftarrow \theta_t + \alpha \nabla_\theta J(\theta_t)$$

- Convergence to local optimum guaranteed
- **Intermediate policies may be arbitrarily bad!**
- **Safe Exploration** [Amodei et al., 2016]: limit risks/costs of exploratory behavior

MOTIVATION

- A working controller is provided as initial policy
- **Fine tuning**: improve it *online* via policy gradient
- Intermediate policies should never be (too much) worse than the initial one

STATE OF THE ART

Existing **safe PG** approaches [Pirota et al., 2013, Papini et al., 2017, 2019]:

- Apply **only** to **Gaussian** policies with **fixed variance**:

$$\pi_\theta(a|s) \sim \mathcal{N}(\mu_\theta(s), \sigma^2)$$

⚠ The variance parameter regulates exploration and has a big impact on convergence speed

- Focus on **monotonic improvement** guarantees:

$$J(\theta_{t+1}) - J(\theta_t) \geq 0$$

⚠ Very strict: **exploration** is totally **sacrificed** due to its immediate costs

CONTRIBUTIONS

- We adopt **a more general definition of safety** that leaves room for exploration
- We propose a **surrogate objective** for the policy variance that takes the **long-term benefits of exploration** into account
- We extend the existing improvement guarantees for Gaussian policies to the **adaptive-variance** case

SETTING

- *Shallow* Gaussian policy parametrization:

$$\pi_\theta(a|s) = \frac{1}{\sqrt{2\pi}\sigma_\omega} \exp \left\{ -\frac{1}{2} \left(\frac{a - \mu_v(s)}{\sigma_\omega} \right)^2 \right\} \quad \mu_v(s) = v^T \phi(s) \quad \sigma_\omega = e^\omega \quad \theta = \begin{cases} v & \text{mean parameter} \\ \omega & \text{variance parameter} \end{cases}$$

- **Safety requirement** (similar to Thomas et al. [2015]):

$$J(\theta_{t+1}) - J(\theta_t) \geq C_t \quad \text{with probability at least } 1 - \delta$$

$$C_t \geq 0 : \quad \text{required improvement}$$

$$C_t < 0 : \quad \text{bounded worsening}$$

- Base algorithm: (normalized) REINFORCE with separate mean and variance updates

$$\begin{cases} v_{t+1} \leftarrow v_t + \alpha_t \nabla_v J(v_t, \omega_t) / \|\nabla_v J(v_t, \omega_t)\| \\ \omega_{t+1} \leftarrow \omega_t + \eta_t \nabla_\omega J(v_{t+1}, \omega_t) / \|\nabla_\omega J(v_{t+1}, \omega_t)\| \end{cases} \quad \text{naive update: too greedy!}$$

- **Adaptive PG**: we look for the *largest* step sizes satisfying the requirement at each iteration

ADAPTIVE EXPLORATION (HEURISTIC)

- We make α_t proportional to σ^2 to exploit its **smoothing effect** [Ahmed et al., 2019]:
- We introduce a **surrogate objective** for ω that accounts for the long-term advantages of exploration:

$$\alpha_t = \alpha \sigma_{\omega_t}^2$$

$$\mathcal{L}(v_t, \omega_t) = \underbrace{J(v_{t+1})}_{\text{performance after next mean update}} \simeq J(v_t, \omega_t) + \alpha \sigma_{\omega_t}^2 \|\nabla_v J(v_t, \omega_t)\|$$

Meta-Exploring Policy Gradient (MEPG):

$$v_{t+1} \leftarrow v_t + \alpha \sigma_{\omega_t}^2 \nabla_v J(v_t, \omega_t) / \|\nabla_v J(v_t, \omega_t)\|, \quad \omega_{t+1} \leftarrow \omega_t + \eta \underbrace{\nabla_\omega \mathcal{L}(v_{t+1}, \omega_t)}_{\text{meta-gradient}} / \|\nabla_\omega \mathcal{L}(v_{t+1}, \omega_t)\|$$

- Learning behavior still depends on **hyperparameters** α, η (step sizes)

SAFE EXPLORATION

- We extend improvement guarantees [Papini et al., 2019] for Gaussian policies to the adaptive-variance setting to obtain adaptive **safe step sizes** for **MEPG**:

$$\begin{cases} \alpha_t = \frac{\sigma_{\omega_t}^2}{F} \left(1 + \sqrt{1 - C_t / C_t^*} \right) \\ \eta_t = \frac{|\lambda_t|}{G} \left(\text{sign}(\lambda_t) + \sqrt{1 - C_t / C_t^*} \right) \end{cases}$$

$$\lambda_t = \frac{\nabla_\omega J(v_t, \omega_t)^T \nabla_\omega \mathcal{L}(v_t, \omega_t)}{\|\nabla_\omega \mathcal{L}(v_t, \omega_t)\|} \quad (\text{scalar projection})$$

- F, G : smoothing constants $\mathcal{O}(R_{\max} \phi_{\max}^2 (1 - \gamma)^{-3})$

- C_t^* : maximum ensurable improvement ($C_t \leq C_t^*$)

- η_t can be negative when exploration constrains with immediate improvement (depending on C_t)

- The resulting **SEPG algorithm (Safely Exploring Policy Gradient)** guarantees $J(v_{t+1}, \omega_{t+1}) - J(v_t, \omega_t) \geq C_t$
- We devise **high-probability** variants for the approximate setting

REFERENCES

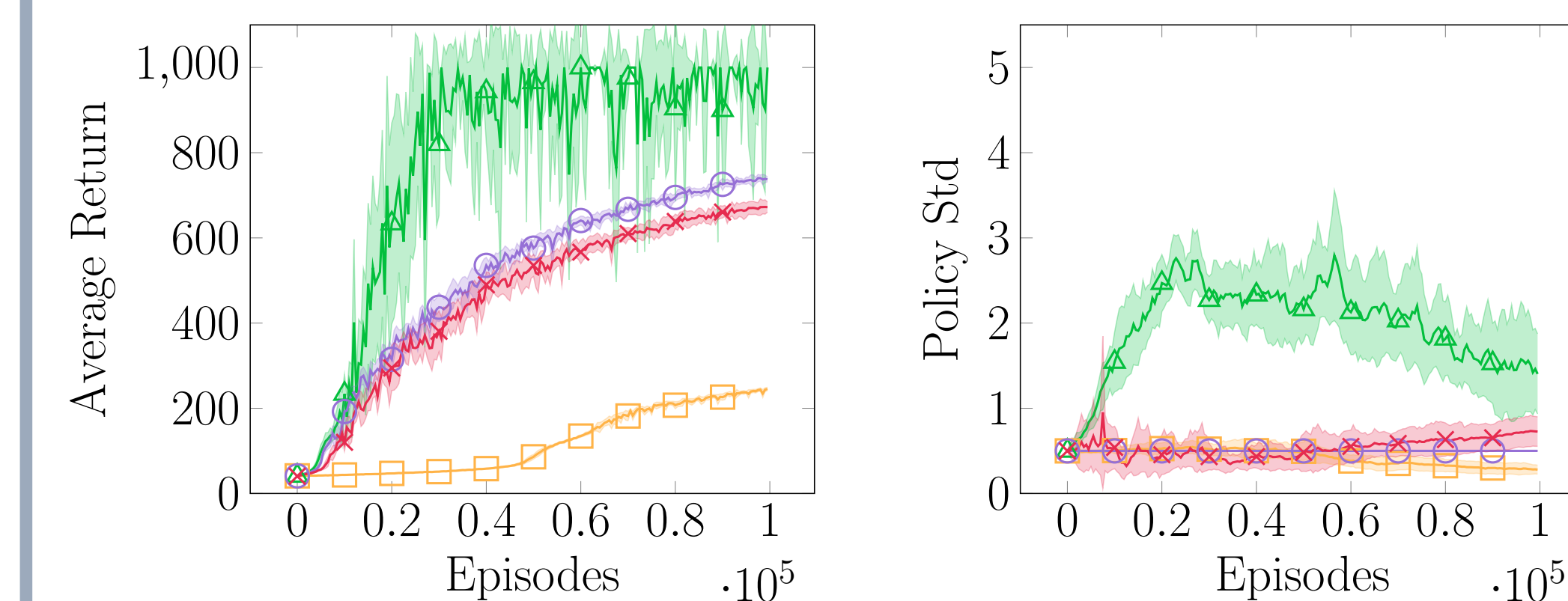
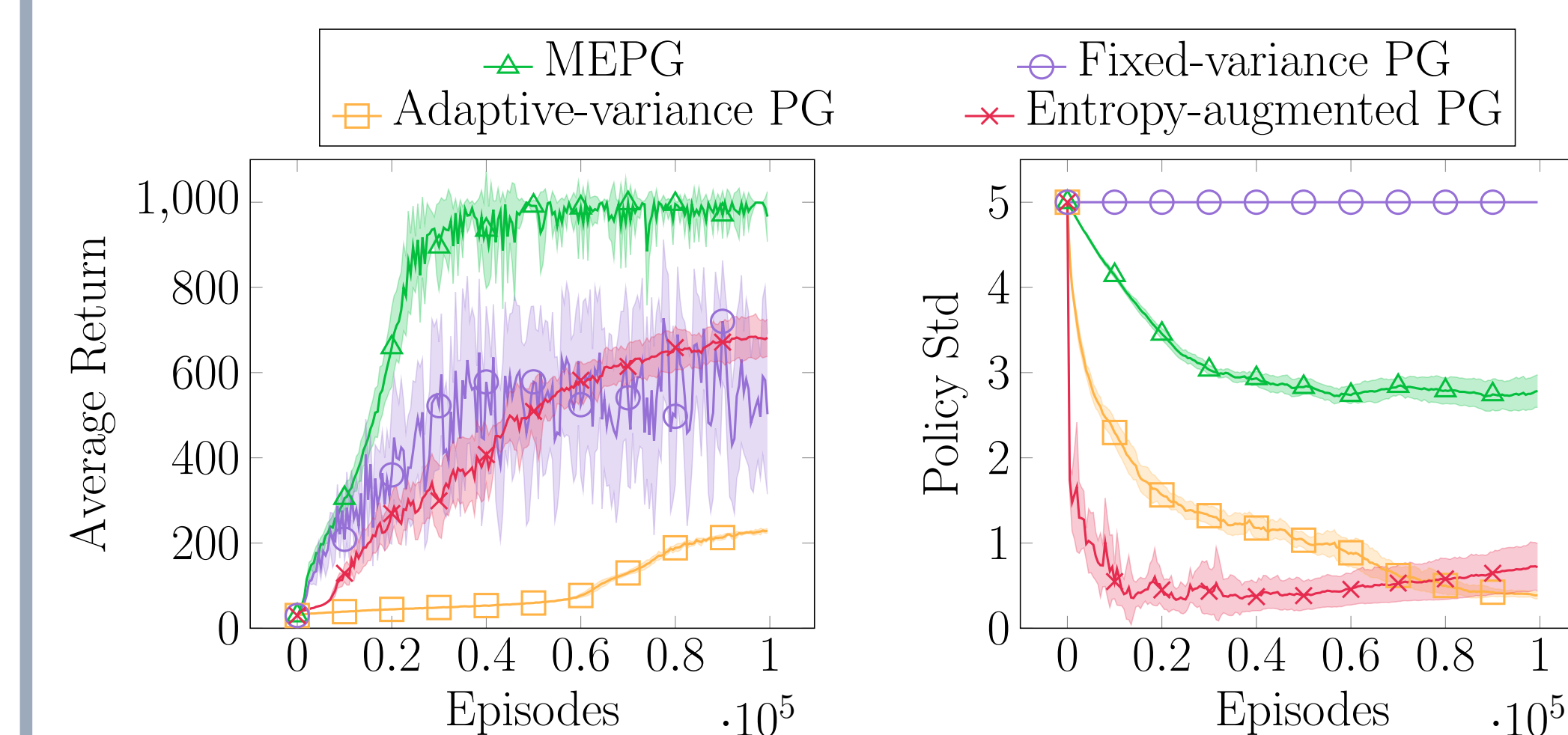
Zafarali Ahmed, Nicolas Le Roux, Mohammad Norouzi, and Dale Schuurmans. Understanding the impact of entropy on policy optimization. In *ICML*, 2019.
Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. 2016.
Marc Peter Deisenroth, Gerhard Neumann, Jan Peters, et al. A survey on policy search for robotics. *Foundations and Trends® in Robotics*, 2(1-2):1-142, 2013.
Matteo Papini, Matteo Pirota, and Marcello Restelli. Adaptive batch size for safe policy gradients. In *NeurIPS*, 2017.
Matteo Papini, Matteo Pirota, and Marcello Restelli. Smoothing policies and safe policy gradients. *arXiv preprint arXiv:1905.03231*, 2019.
Matteo Pirota, Marcello Restelli, and Luca Bascetta. Adaptive step-size for policy gradient methods. In *NeurIPS*, pages 1394-1402, 2013.
Philip Thomas, Georgios Theodorou, and Mohammad Ghavamzadeh. High confidence policy improvement. In *ICML*, 2015.

SPECIAL REQUIREMENTS

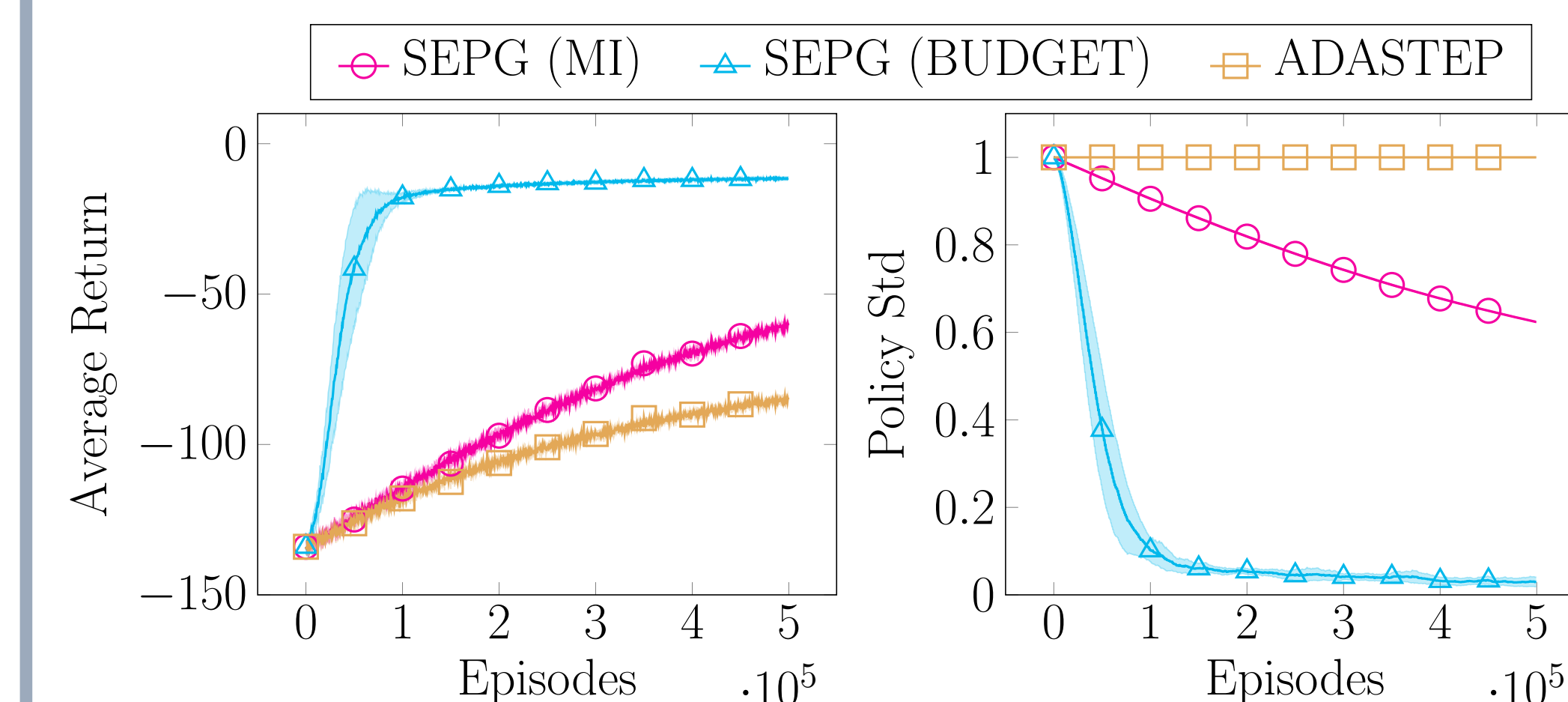
- **Monotonic improvement (MI)**: $C_t \equiv 0$
- **Fixed improvement**: constant C_t (possibly < 0)
- **Fixed threshold**: $C_t = J_{\min} - J(\theta_t)$
- **Fine tuning (BUDGET)**: $C_t = J(\theta_0) - J(\theta_t)$

EXPERIMENTS

MEPG on Cart-Pole



SEPG on LQ



SEPG on Cart-Pole

