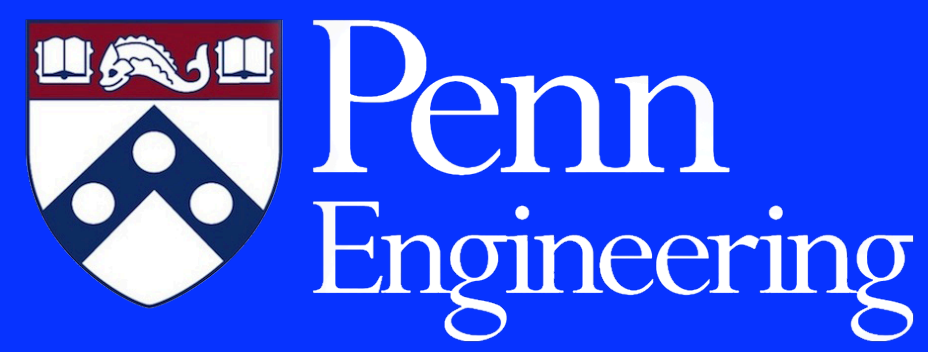


Safe Policy Search for Lifelong Reinforcement Learning with Sublinear Regret



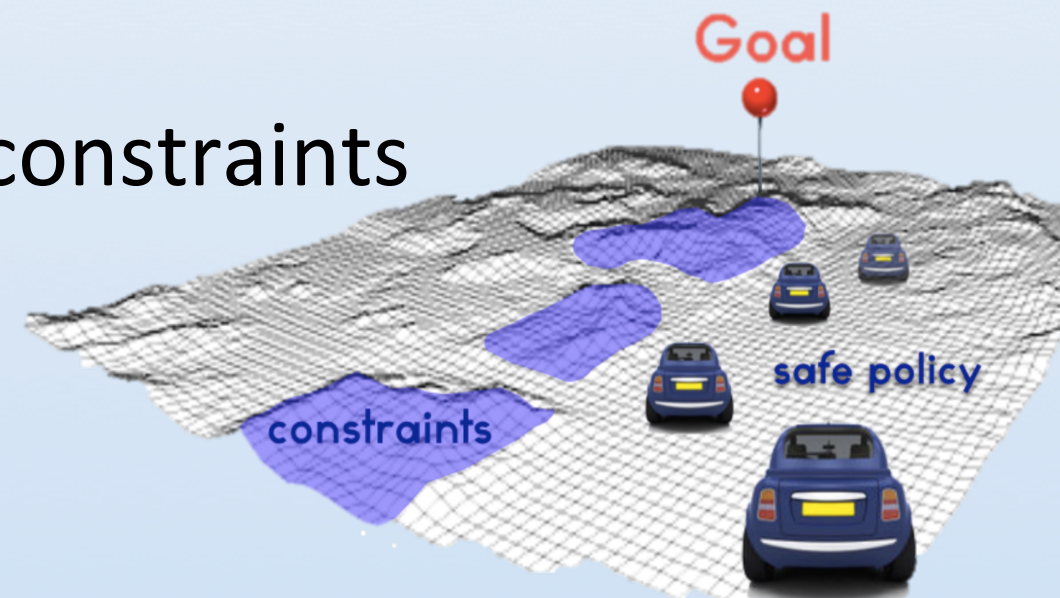
Summary

We developed a lifelong policy gradient learner that operates in an adversarial setting to learn multiple tasks online while enforcing safety constraints on the learned policies.

- Fully online learning of multiple, consecutive RL tasks
- Ensures “safe” policies by respecting safety constraints
- Exhibits *sublinear regret* for lifelong policy search
- Validated on benchmark dynamical systems and quadrotor control

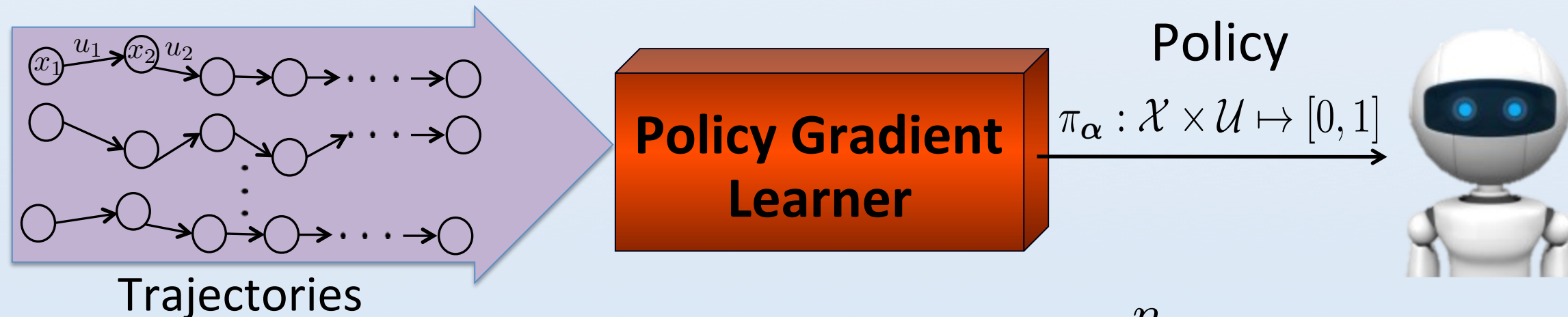
Motivations

1. Reuse knowledge from previously learned tasks to accelerate the learning of new control policies
 - Lifelong RL to learn multiple, consecutive tasks online
 - Exhibit vanishing regrets
2. Robotic control policies must obey safety constraints (e.g., prevent damage to the robot and environment, avoid catastrophic failure)
 - Incorporate constraints directly into the optimization



Background: Policy Gradient (PG) Methods

- Agent interacts with environment, taking consecutive actions
- PG methods support continuous state and action spaces
 - Have shown recent success in applications to robotic control



Goal: find policy π_α that minimizes $l(\alpha) = \sum_{k=1}^n p_\alpha(\tau^{(k)}) C(\tau^{(k)})$

$$p_\alpha(\tau^{(k)}) = \mathcal{P}_0(x_0^{(k)}) \prod_{m=0}^{M-1} \mathcal{P}(x_{m+1}^{(k)} | x_m^{(k)}, u_m^{(k)}) \pi_\alpha(u_m^{(k)} | x_m^{(k)})$$

$$C(\tau^{(k)}) = \frac{1}{M} \sum_{m=0}^{M-1} c_{m+1}^{(k)}$$

Background: Online Learning & Regret Analysis

Regret Minimization Game: Each round $j = 1 \dots R$,

- agent chooses a prediction θ_j , and
- environment (i.e., the adversary) chooses a loss function l_j

Goal: minimize cumulative regret (modified for multi-task case)

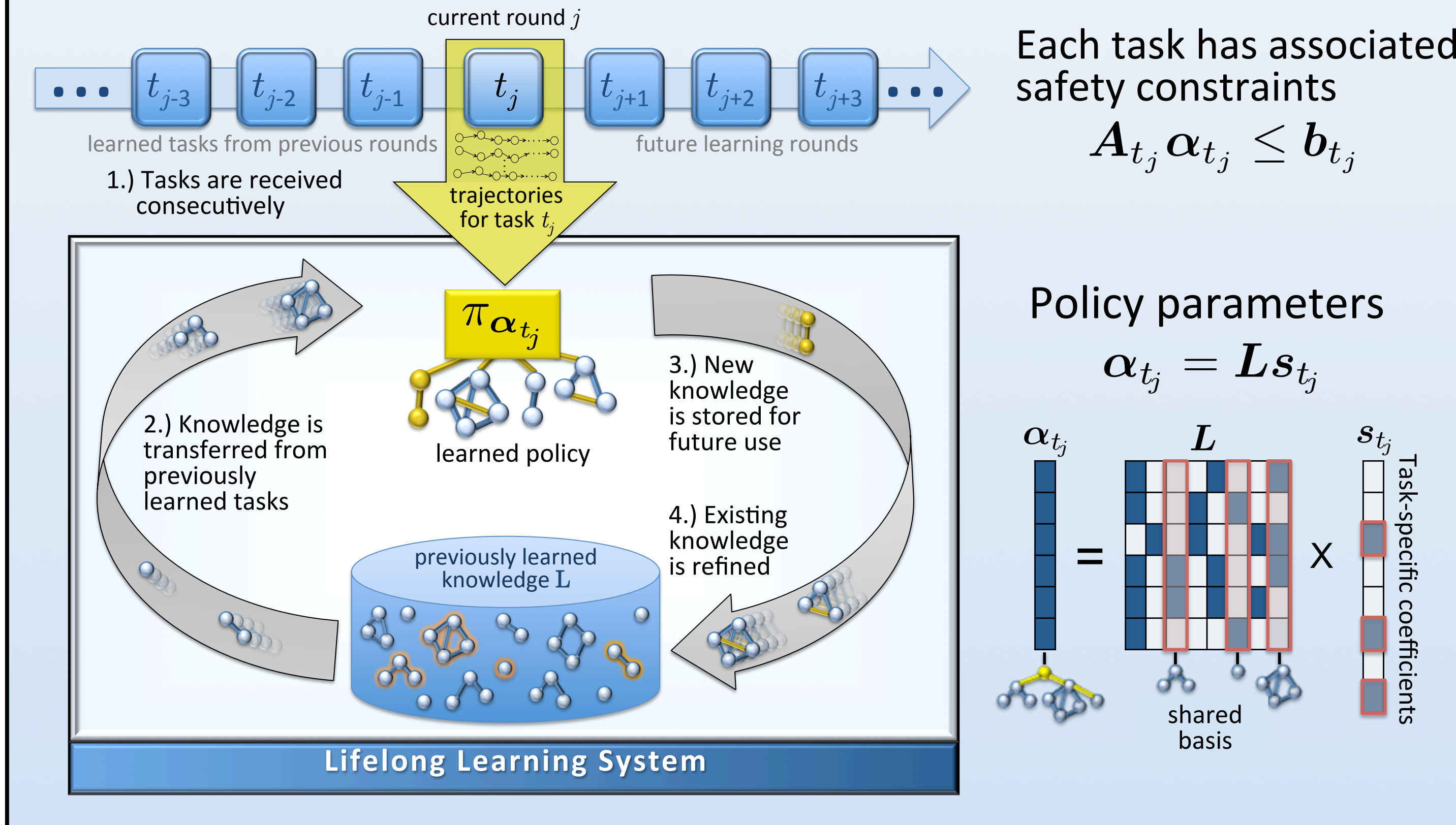
$$\mathfrak{R}_R = \sum_{j=1}^R l_{t_j}(\theta_j) - \inf_{u \in \mathcal{K}} \left[\sum_{j=1}^R l_{t_j}(u) \right]$$

l_{t_j} : loss of task t at round j

Solve via “Follow the Regularized Leader”:

- 1.) Find $\tilde{\theta}$ via unconstrained optimization over accumulated losses
- 2.) Project $\tilde{\theta}$ onto the constraint set via Bregman projections

Lifelong Learning Framework



Safe Lifelong Policy Search

Multi-task Optimization Problem after observing r rounds:

$$\min_{L, S} \sum_{j=1}^r [\eta_{t_j} l_{t_j}(L s_{t_j})] + \mu_1 \|S\|_F^2 + \mu_2 \|L\|_F^2$$

$$\text{s.t. } A_{t_j} \alpha_{t_j} \leq b_{t_j} \quad \forall t_j \in \mathcal{I}_r$$

$$\lambda_{\min}(L L^T) \geq p \quad \text{and} \quad \lambda_{\max}(L L^T) \leq q$$

Online Formulation

- Let $\theta = [\text{vec}(L) \text{vec}(S)]^T$ be the vector of all parameters
- The MTL objective can be written online as

$$\theta_{r+1} = \arg \min_{\theta \in \mathcal{K}} \Omega_r(\theta) \quad \Omega_0(\theta) = \mu_2 \sum_{i=1}^{dk} \theta_i^2 + \mu_1 \sum_{i=dk+1}^{dk+k|\mathcal{T}|} \theta_i^2$$

$$\Omega_j(\theta) = \Omega_{j-1}(\theta) + \eta_{t_j} l_{t_j}(\theta)$$

where the loss for task t_j is the following bilinear product in θ :

$$l_{t_j}(\theta) = -\frac{1}{n_{t_j}} \sum_{k=1}^{n_{t_j}} \sum_{m=0}^{M_{t_j}-1} \log \left[\pi_{\Theta_L \Theta_{s_{t_j}}}^{(t_j)}(u_m^{(k), t_j} | x_m^{(k), t_j}) \right]$$

$$\Theta_L = \begin{bmatrix} \theta_1 & \dots & \theta_{d(k-1)+1} \\ \vdots & \vdots & \vdots \\ \theta_d & \dots & \theta_{dk} \end{bmatrix}, \quad \Theta_{s_{t_j}} = \begin{bmatrix} \theta_{dk+1} \\ \vdots \\ \theta_{(d+1)k+1} \end{bmatrix}$$

Online Solution

- Step 1: Unconstrained policy optimization via alternating optimization over L and S to yield unconstrained solution $\tilde{\theta}_{r+1}$
- Step 2: Project $\tilde{\theta}_{r+1}$ to the constraint set via the Bregman divergence over Ω_r — involves solving 2nd order cone and semi-definite programs

Main Theoretical Result

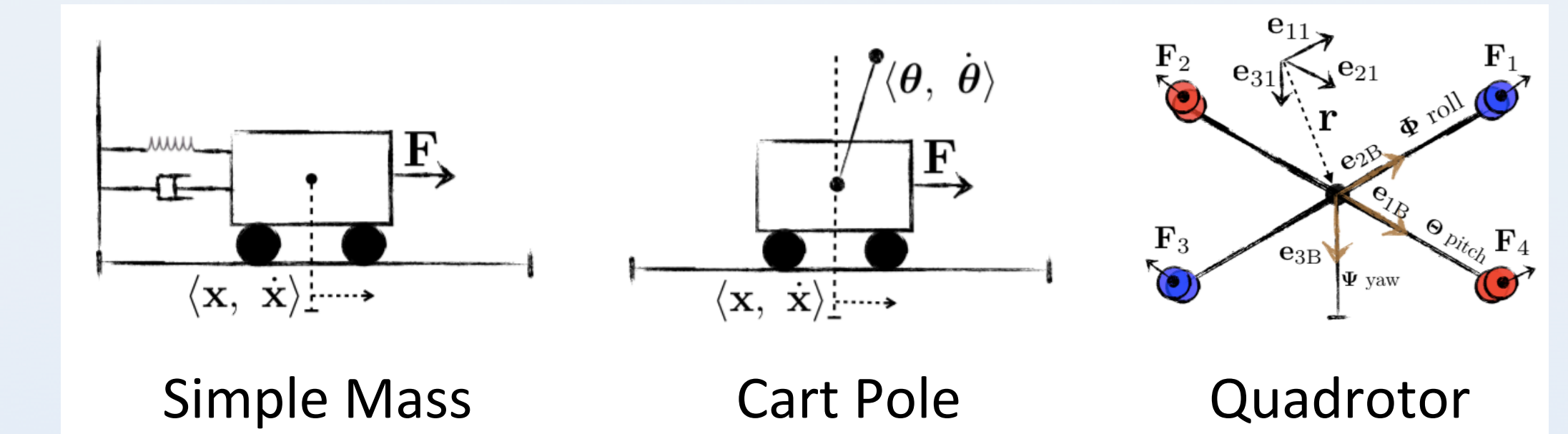
We prove that the safe lifelong policy search algorithm has sublinear regret of $\mathcal{O}(\sqrt{R})$ in the total number of rounds R .

Theorem 1 (Sublinear Regret). After R rounds and choosing $\forall t_j \in \mathcal{I}_R \eta_{t_j} = \eta = \frac{1}{\sqrt{R}}$, $L|_{\theta_1} = \text{diag}_k(\zeta)$, with $\text{diag}_k(\cdot)$ being a diagonal matrix among the k columns of L , $p \leq \zeta^2 \leq q$, and $S|_{\theta_1} = \mathbf{0}_{k \times |\mathcal{T}|}$, the safe lifelong reinforcement learner exhibits sublinear regret of the form:

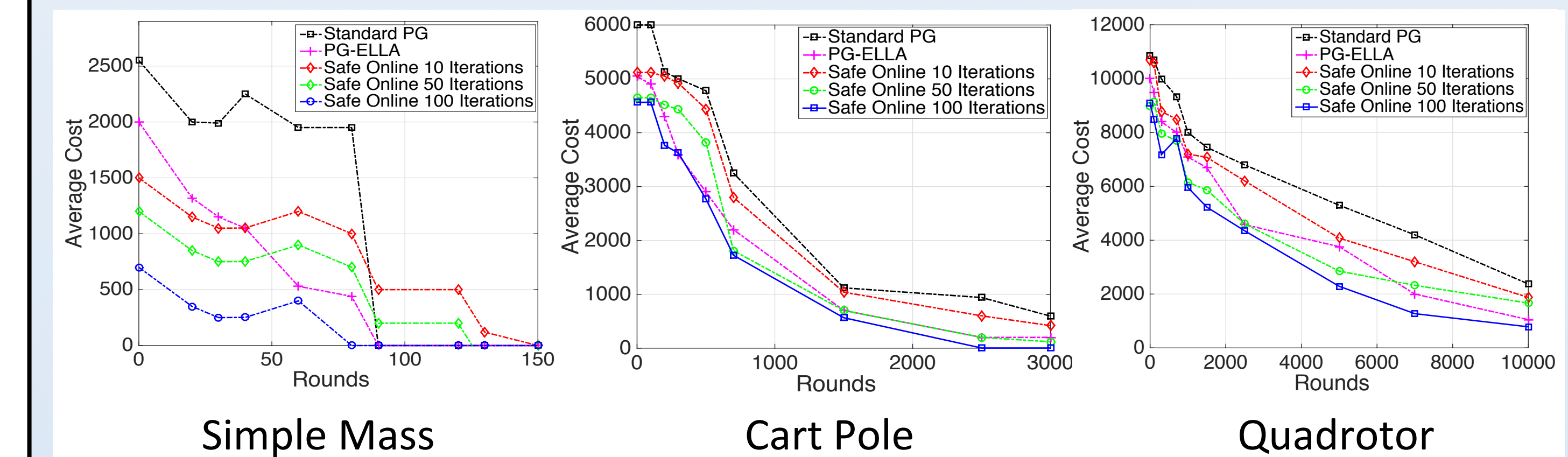
$$\sum_{j=1}^R l_{t_j}(\hat{\theta}_j) - l_{t_j}(u) = \mathcal{O}(\sqrt{R}) \quad \text{for any } u \in \mathcal{K}.$$

Experimental Results

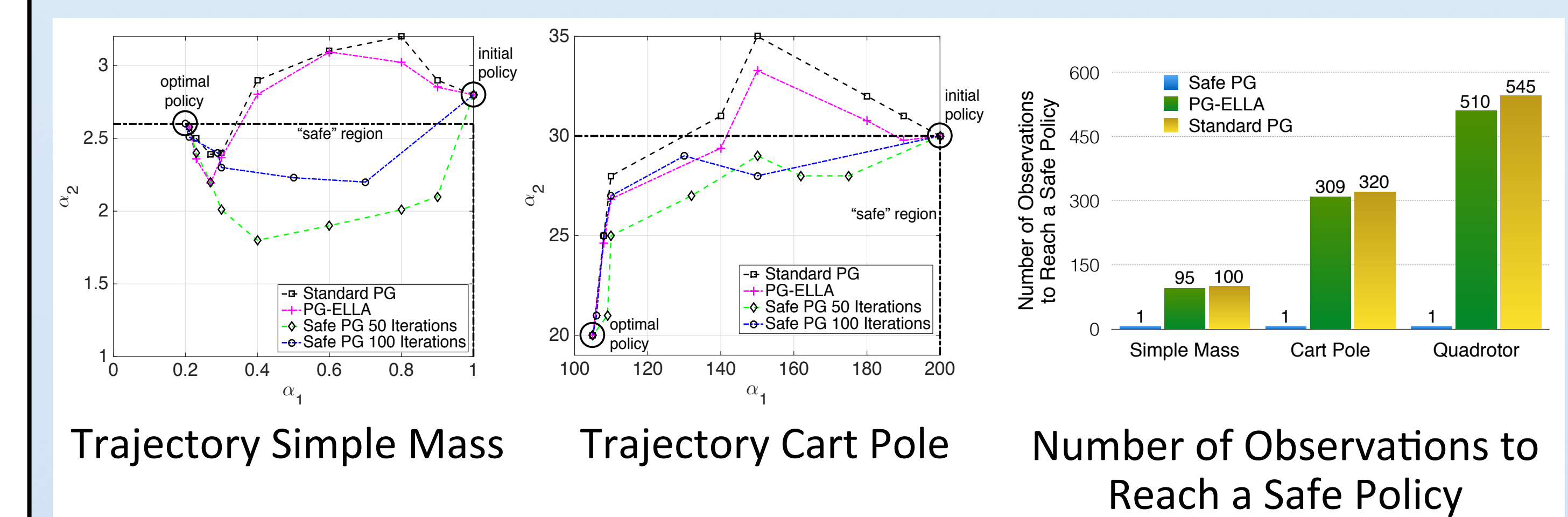
Learn policies for consecutive control tasks on three types of systems



Superior Performance over standard PG and the lifelong learner PG-ELLA



Enforces the Given Safety Constraints, unlike alternative methods



- Note that our approach immediately projects policies to safe regions even during the policy search process, unlike other methods

Acknowledgements

This research was supported by ONR grant #N00014-11-1-0139 and AFRL grant #FA8750-14-1-0069.

We thank Ali Jadbabaie for assistance with the optimization solution, and the anonymous reviewers for their helpful feedback.