

Constrained Policy Optimization via Bayesian World Models

Yarden As, Ilnura Usmanova, Sebastian Curi, Andreas Krause



ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

AUTOMATIC
CONTROL
LABORATORY **IFA**



Learning &
Adaptive Systems

Reinforcement learning agents demonstrate high potential in solving complex tasks.

How can we make them *safe*?



Constrained Markov decision processes

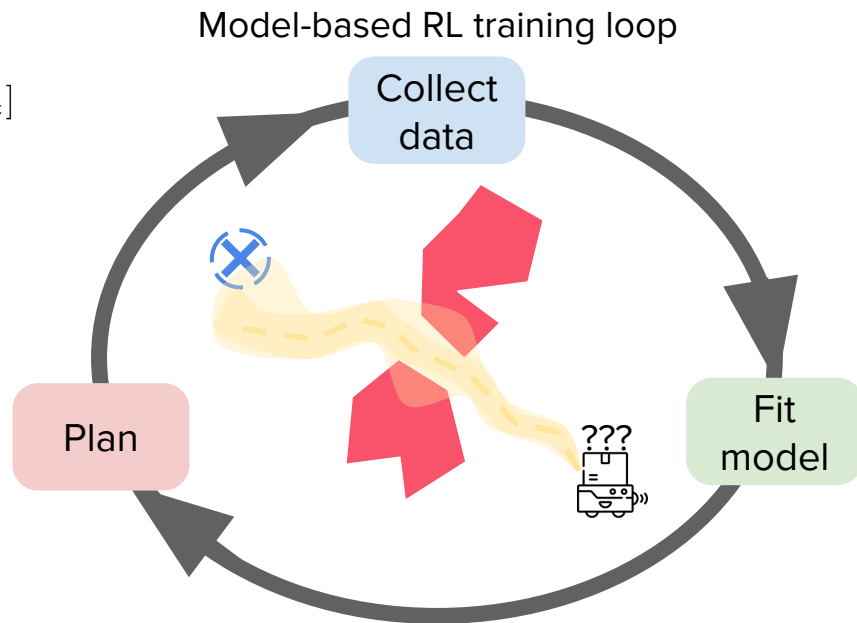
Formulating safety in reinforcement learning

$$\textbf{Goal: } \max_{\pi \in \Pi} J(\pi, p^*) \quad J = \mathbb{E}_{\pi, p^*} [\sum_t r_t]$$

s.t.

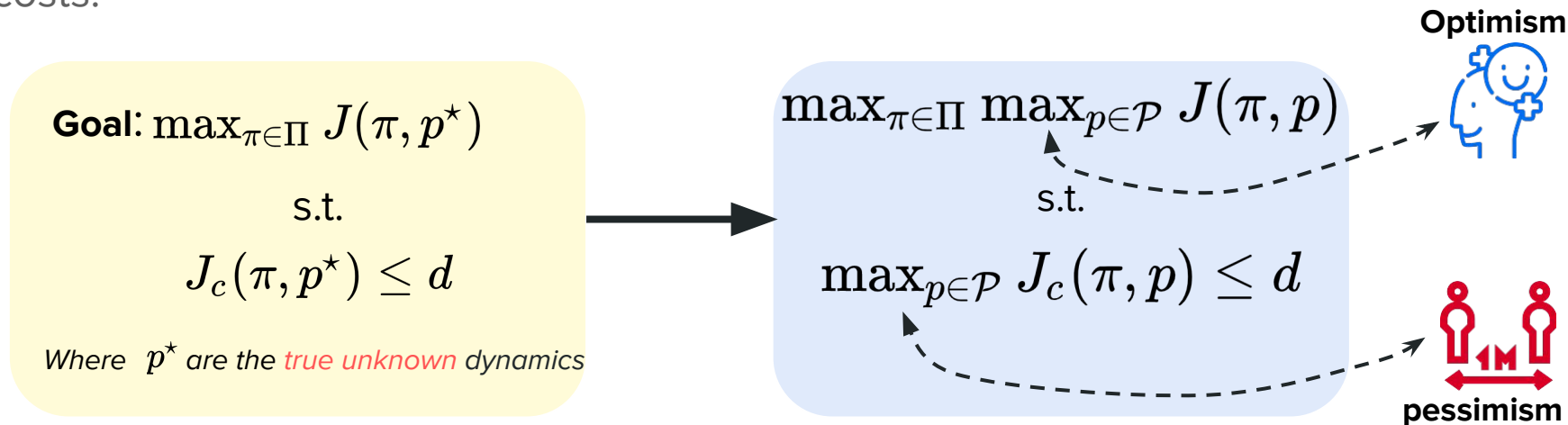
$$J_c(\pi, p^*) \leq d \quad J_c = \mathbb{E}_{\pi, p^*} [\sum_t c_t]$$

Where p^* are the *true unknown* dynamics



How do we use optimism and pessimism?

Idea: use the dynamics to be optimistic for the rewards and pessimism for the costs.



But the inner maximization is intractable...



Upper bounds estimation via posterior sampling

$$\theta \sim p(\theta|\mathcal{D})$$



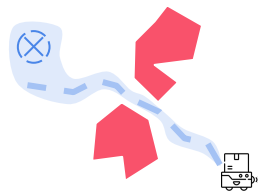
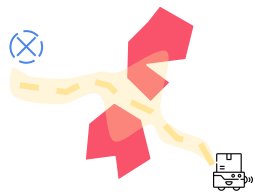
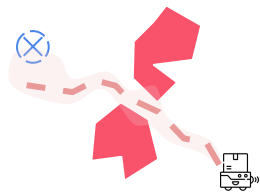
$$p(s_{t+1}|s_t, a_t, \theta)$$



$$p(s_{t+1}|s_t, a_t, \theta)$$



$$p(s_{t+1}|s_t, a_t, \theta)$$



Cost



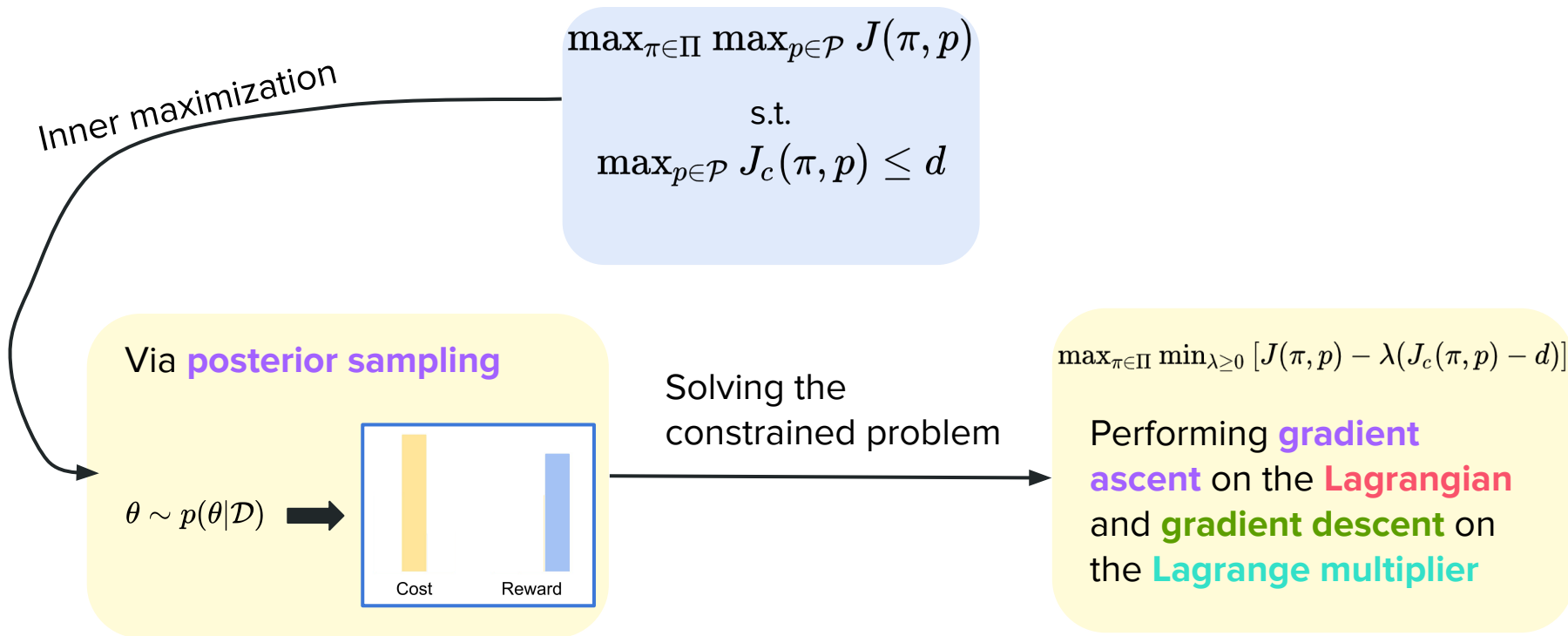
Reward



1. Sample model.
2. Sample trajectories.
3. Evaluate. (via $J(\pi, p) \approx \mathbb{E}_{\pi, p} \left[\frac{1}{H} \sum_t^H V(s_t) \right]$
 $J_c(\pi, p) \approx \mathbb{E}_{\pi, p} \left[\frac{1}{H} \sum_t^H V_c(s_t) \right]$).
4. Repeat steps 1-3.
5. Obtain optimistic/pessimistic evaluations.

SWAG: Maddox et al. (2019). A Simple Baseline for Bayesian Uncertainty in Deep Learning. **RSSM:** Hafner et al. (2019). Learning Latent Dynamics for Planning from Pixels. ICML (2019)

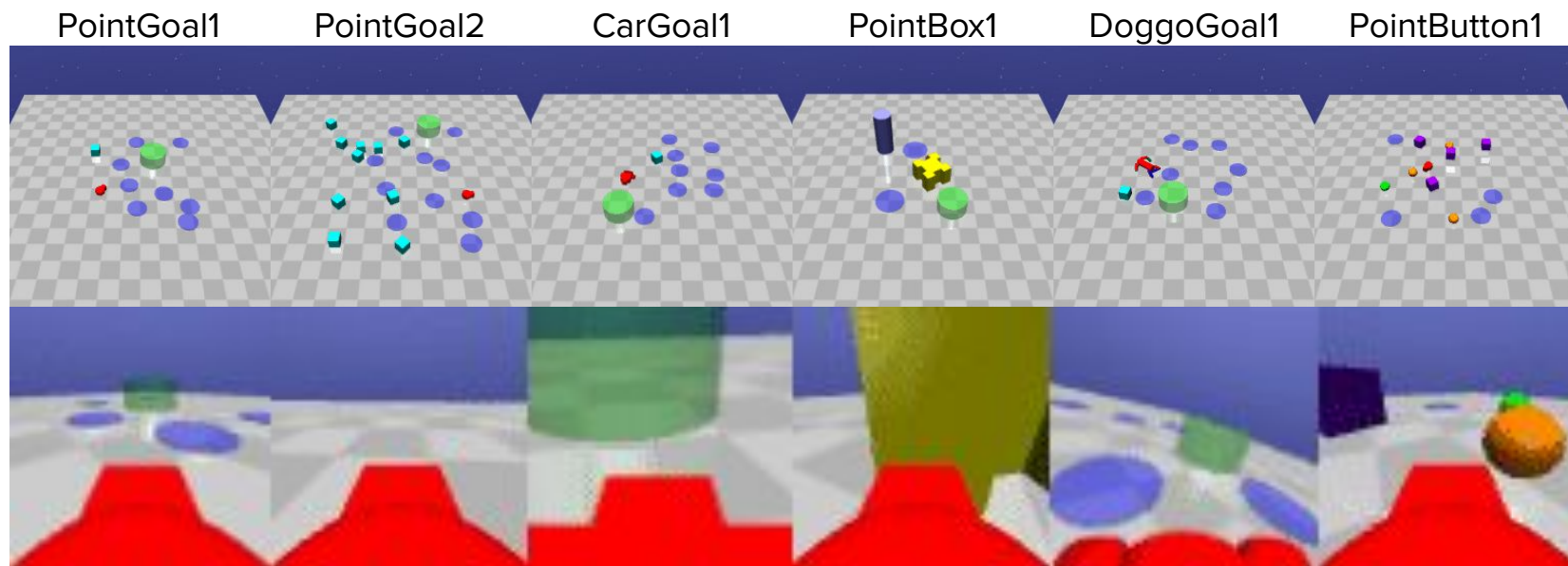
Lagrangian Model-Based Agent (LAMBDA)



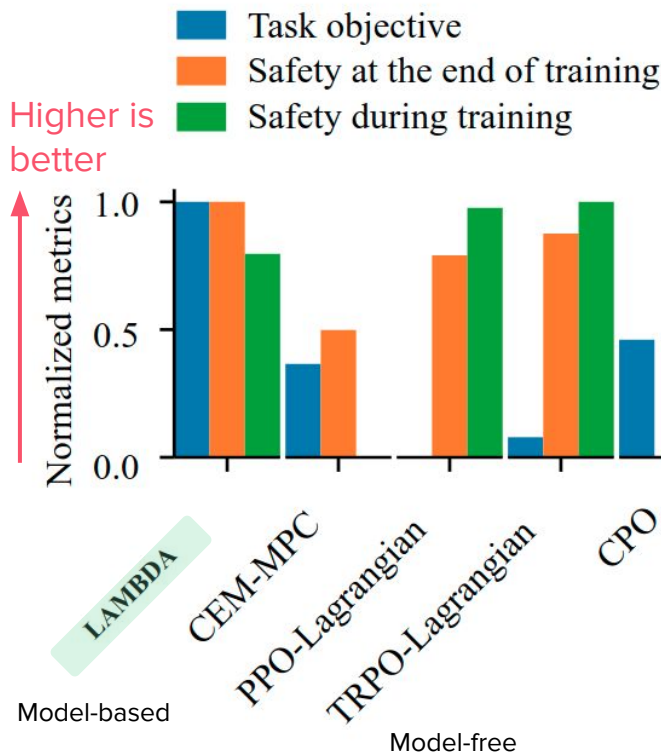
Nocedal & Wright (2006). *Numerical Optimization*. Springer, 2nd Edition.

Benchmark with Safety-Gym

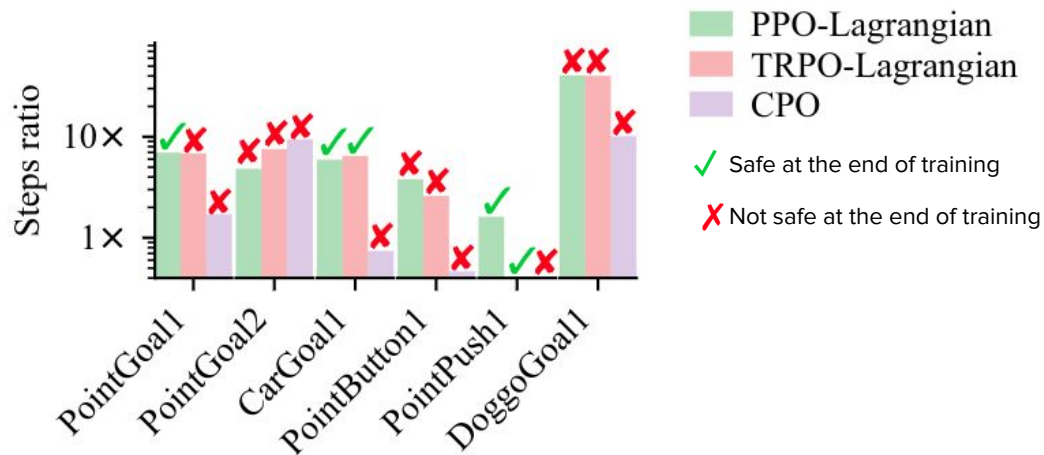
Robots, tasks and observations



Experimental results



Required amount of steps to reach LAMBDA's performance

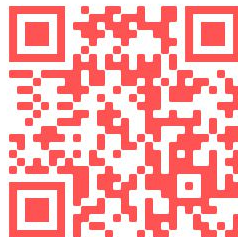


Between **x5-10** more sample efficient

Same amount of training

LAMBDA can *efficiently* learn to solve *complex* tasks which require *safe* behavior, from image observations.

Paper (pre-print, arXiv):



Code (github):

