

Trustworthy AI Autonomy

M3-1: Safe Reinforcement Learning

Ding Zhao

Assistant Professor
Carnegie Mellon University

Plan for today

- Safe RL methods
 - Safety Gym environment
 - Trust-region search: TRPO, PPO
 - Constrained MDP

OpenAI Safety Gym

Benchmarking Safe Exploration in Deep Reinforcement Learning

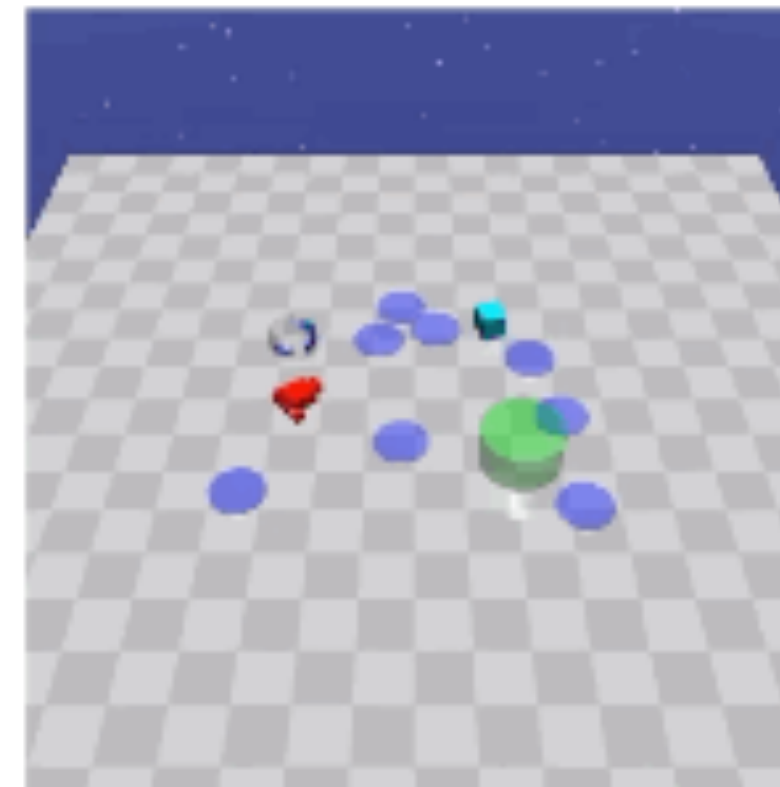
Alex Ray*
OpenAI

Joshua Achiam*
OpenAI

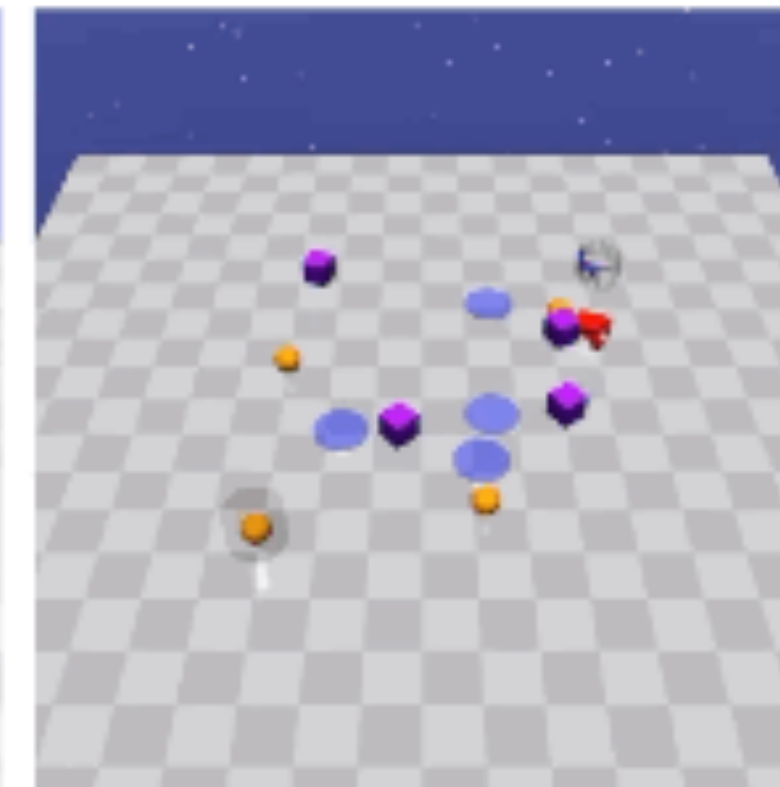
Dario Amodei
OpenAI

Point is a simple robot constrained to the 2D plane, with one actuator for turning and another for moving forward or backward. Point has a front-facing small square which helps with the Push task.

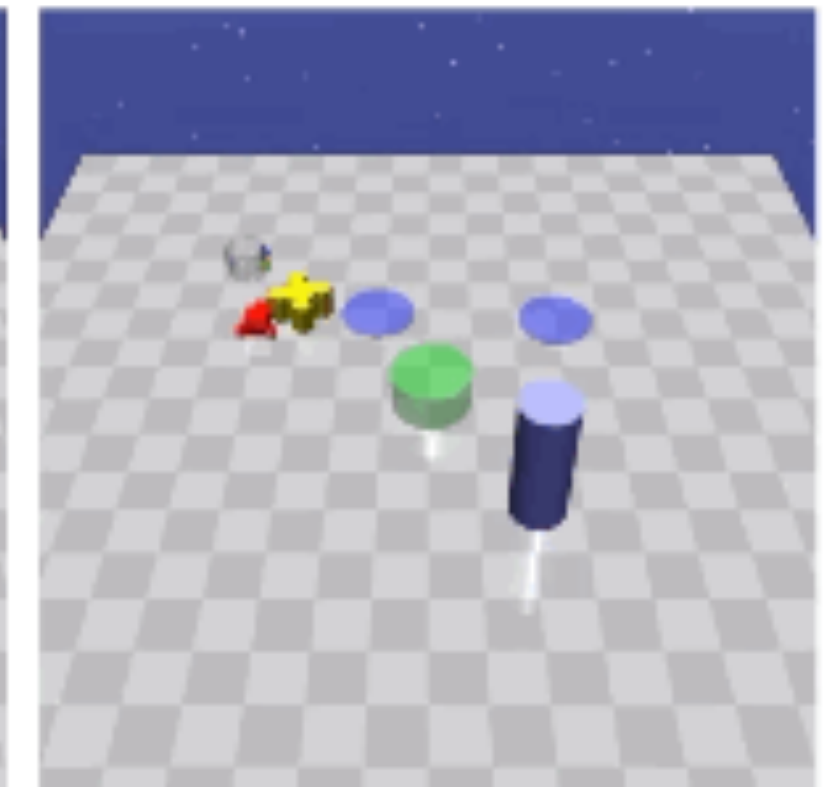
Car has two independently-driven parallel wheels and a free-rolling rear wheel. For this robot, turning and moving forward or backward require coordinating both of the actuators.



Goal: Move to a series of goal positions.

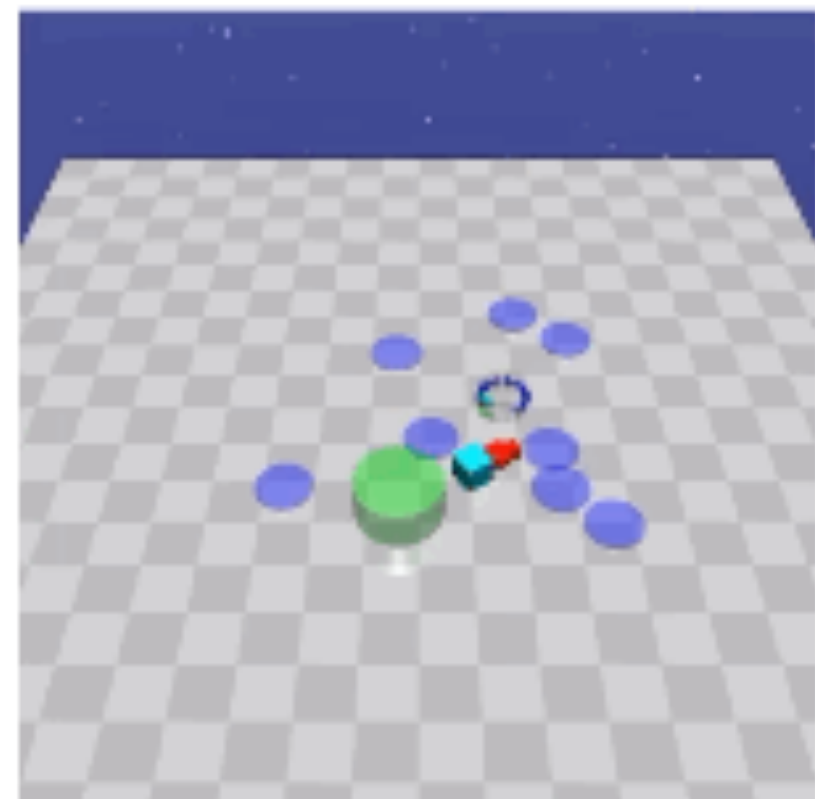


Button: Press a series of goal buttons.

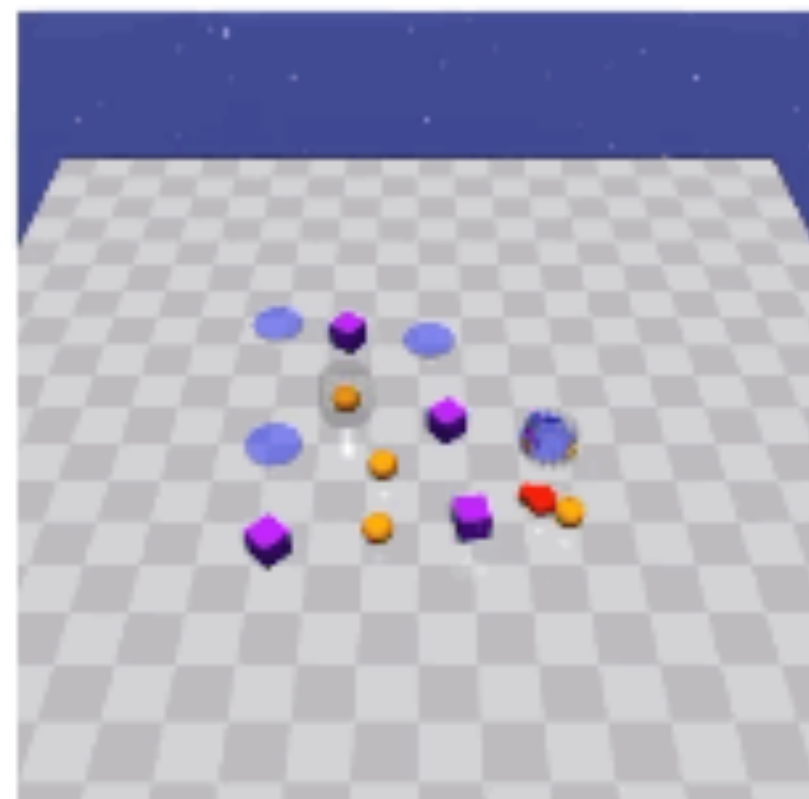


Push: Move a box to a series of goal positions.

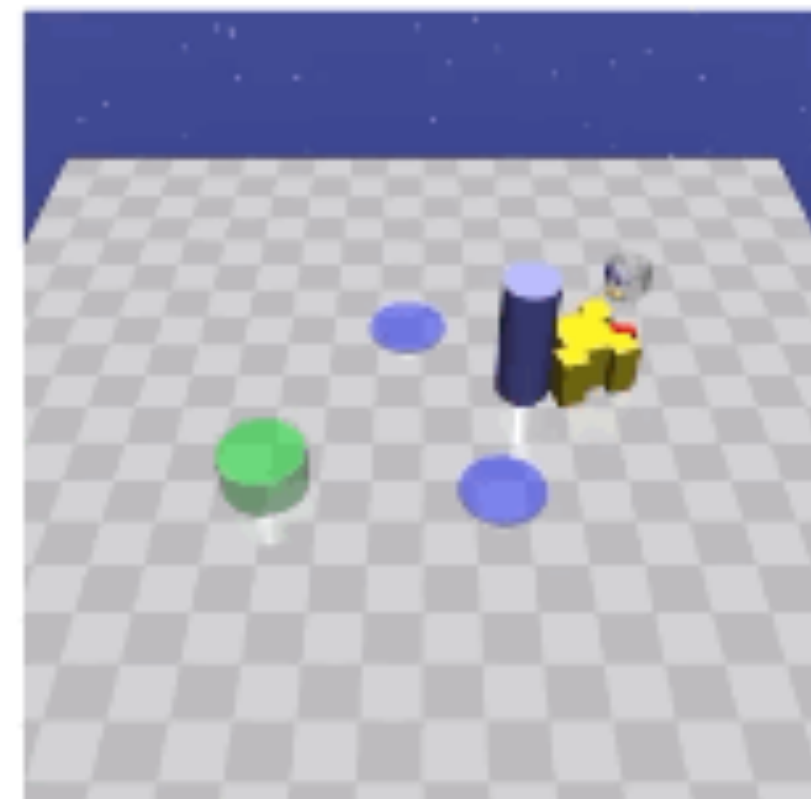
Doggo is a quadruped with bilateral symmetry. Each of its four legs has two controls at the hip, for azimuth and elevation relative to the torso, and one in the knee, controlling angle. A uniform random policy keeps the robot from falling over and generates travel.



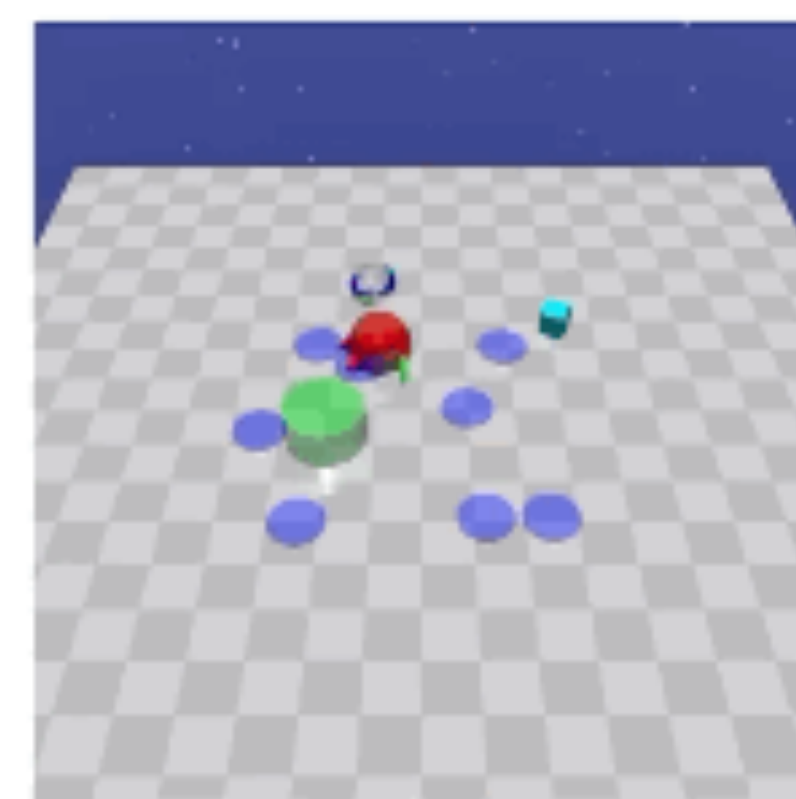
Goal: Move to a series of goal positions.



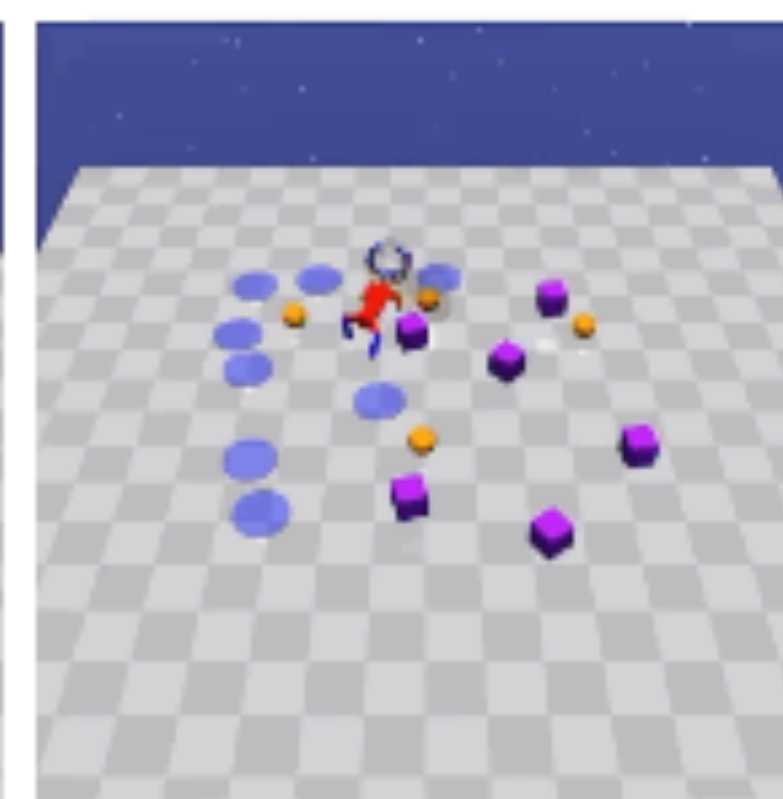
Button: Press a series of goal buttons.



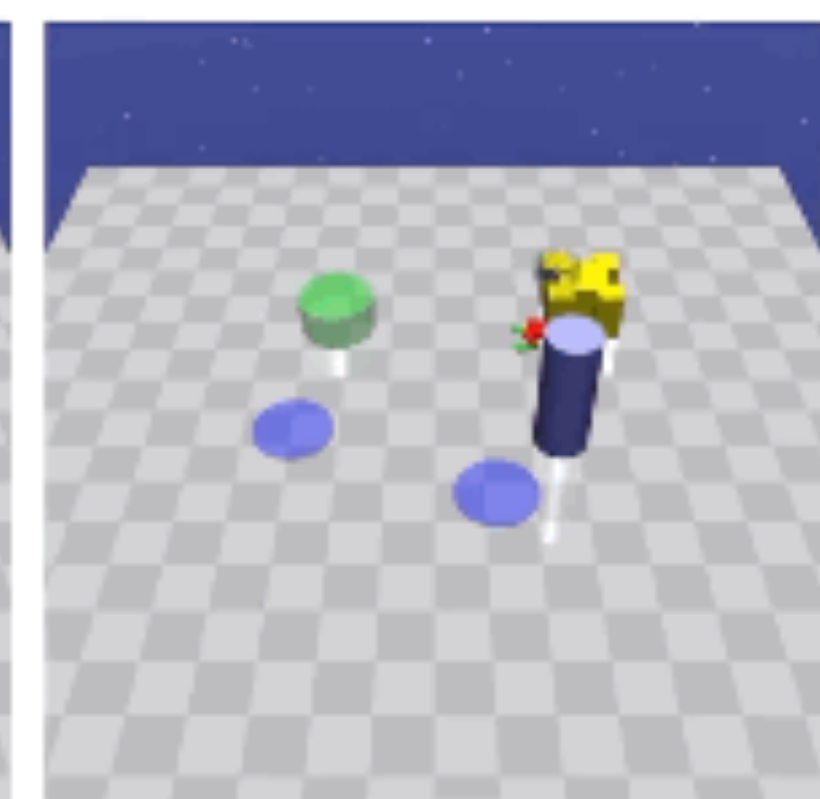
Push: Move a box to a series of goal positions.



Goal: Move to a series of goal positions.



Button: Press a series of goal buttons.



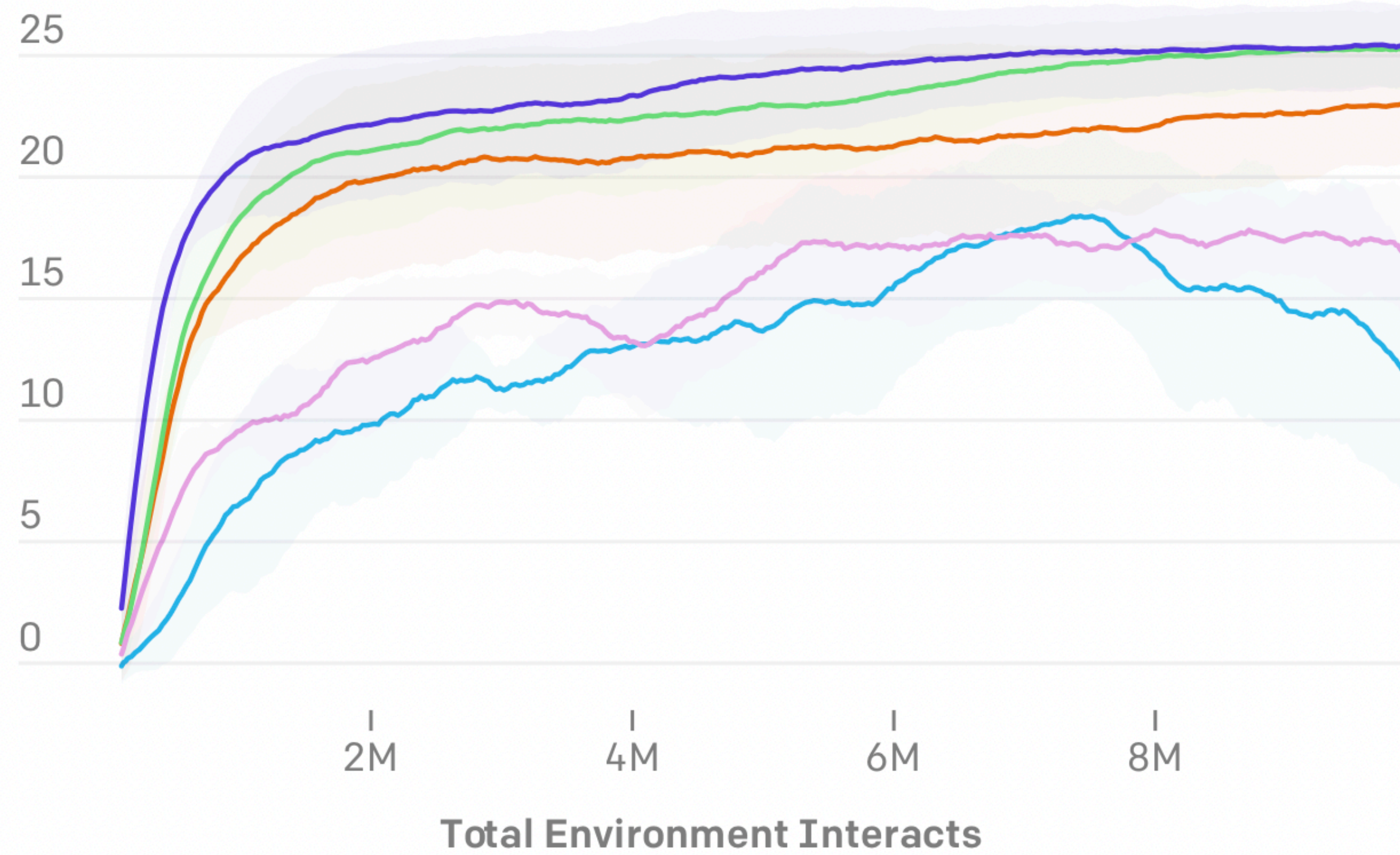
Push: Move a box to a series of goal positions.

Return and cost trade off against each other meaningfully

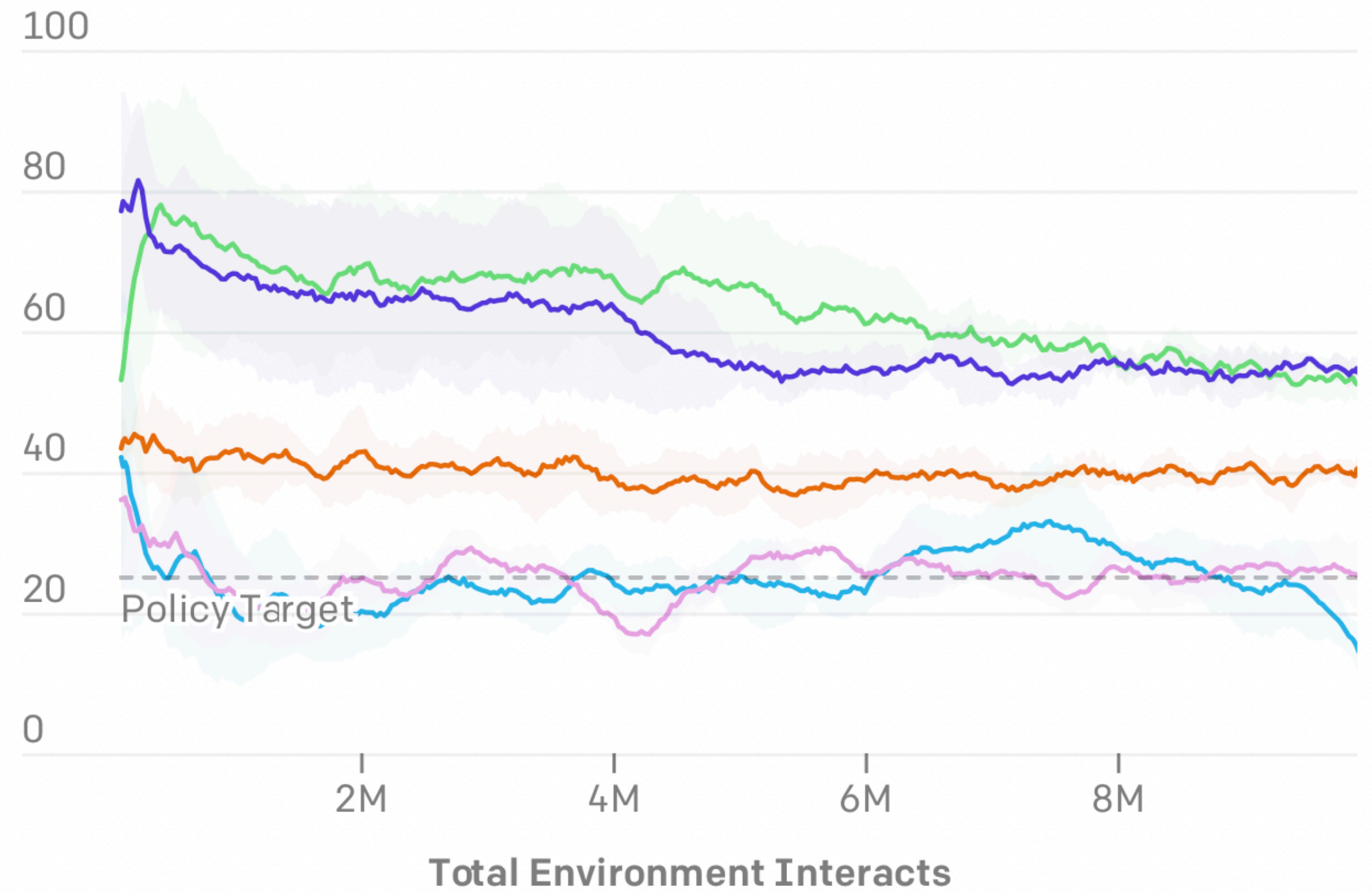
Robot Point Task Goal Level 1

● CPO ● PPO ● PPO-Lagrangian ● TRPO ● TRPO-Lagrangian

Return



Cost



Safe Reinforcement Learning

- *“Safe Reinforcement Learning can be defined as the process of learning policies that maximize the expectation of the return in problems in which it is important to ensure reasonable system performance and/or respect safety constraints during the learning and/or deployment processes.”*
- Where to add the constraints?
 - Add constraints to the exploration procedure: TRPO, PPO
 - Add constraints to the optimization criteria (performance/constraints): CMDP

Constrained Markov Decision Process

- Augment the MDP with a set C of auxiliary cost functions, C_1, \dots, C_m (with each one a function $C_i : S \times A \times S \rightarrow \mathbb{R}$ mapping transition tuples to costs, like the usual reward), and limits d_1, \dots, d_m . Let $J_{C_i}(\pi)$ denote the expected discounted return of policy π with respect to cost function

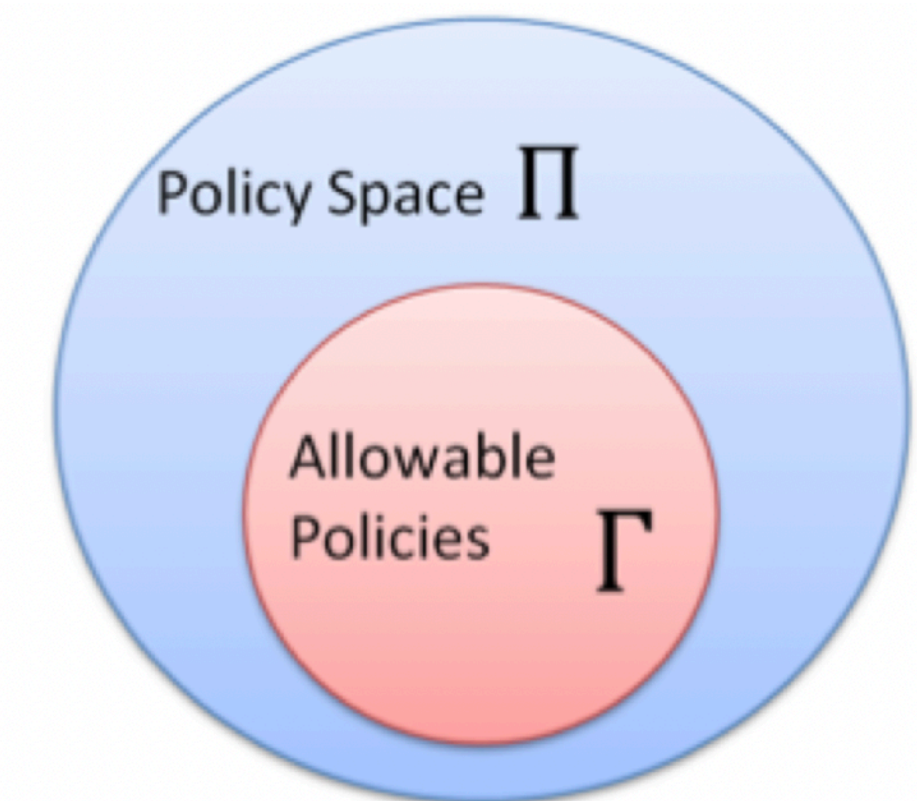
$$C_i : J_{C_i}(\pi) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t C_i(s_t, a_t, s_{t+1}) \right].$$

- The set of feasible stationary policies for a CMDP is then

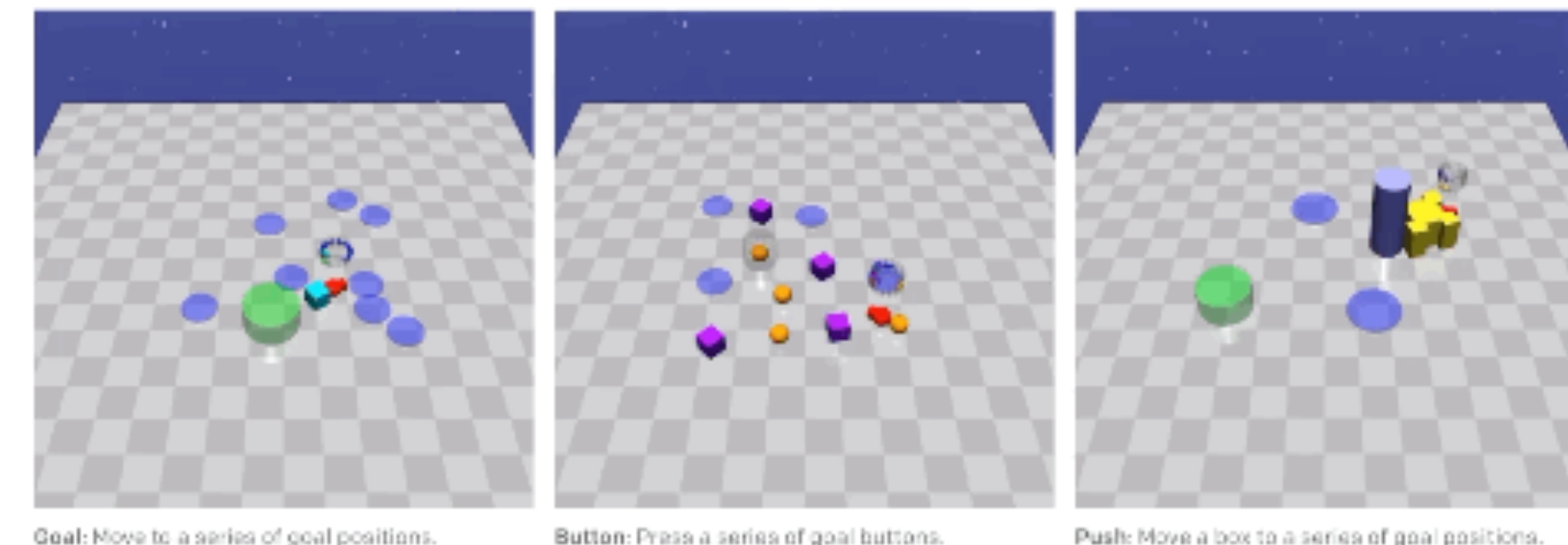
$$\Pi_C \doteq \left\{ \pi \in \Pi : \forall i, J_{C_i}(\pi) \leq d_i \right\}$$

- and the reinforcement learning problem in a CMDP is

$$\pi^* = \arg \max_{\pi \in \Pi_C} J(\pi)$$



Point is a simple robot constrained to the 2D plane, with one actuator for turning and another for moving forward or backward. Point has a front-facing small square which helps with the Push task.



Safe RL methods

- Constrained Markov Decision Process
- Lagrangian TRPO and PPO
 - Trust Region Policy Optimization (TRPO)
 - Proximal Policy Optimization (PPO)
- Constrained Optimization-based Methods

I will focus on the key motivations and intuitions. For details of implementation, you can read these blogs and corresponding papers.

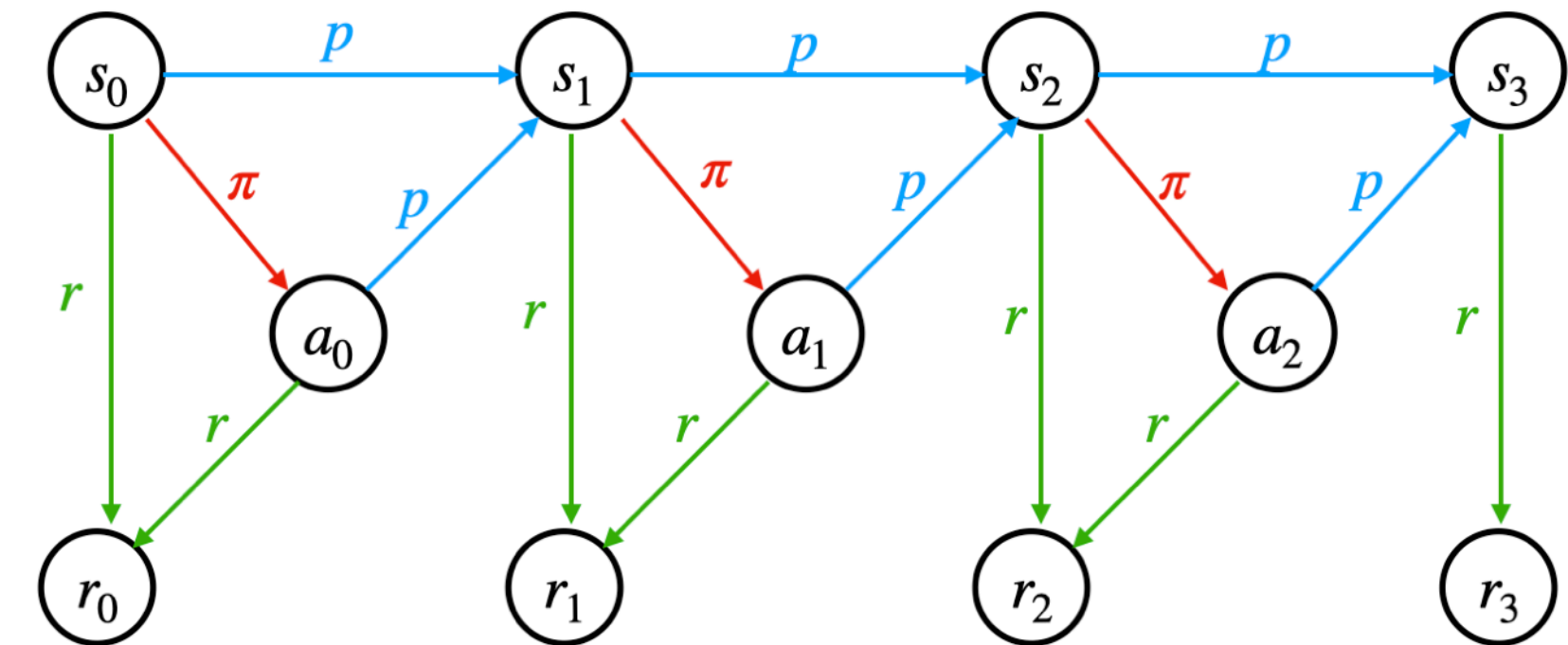
https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html

Recap: Markov Decision Process

- Defined by: $(\mathcal{S}, \mathcal{A}, r, p)$
 - \mathcal{S} : set of possible states
 - \mathcal{A} : set of possible actions
 - r : reward function
 - p : dynamics function
- Goal of MDP: given $(\mathcal{S}, \mathcal{A}, r, p)$, $\rho_0(\cdot)$, T or γ , we want

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

$$s_{t+1} \sim p(\cdot | s_t, a_t)$$
$$a_t \sim \pi(\cdot | s_t)$$
$$r_t \sim r(\cdot | s_t, a_t)$$



Baselines for Safe RL algorithms

- Trust Region Policy Optimization (TRPO)
- Proximal Policy Optimization (PPO)

I will focus on the key motivations and intuitions. For details of implementation, you can read these blogs and corresponding papers.

https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html

Trust Region Policy Optimization (TRPO)

- Safety issue with Vanilla Policy Gradient (VPG)
- $J(\theta, \mathcal{D}_{\pi_\theta}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \pi_\theta \right]$, $\theta^* = \arg \max_{\theta} J(\theta, \mathcal{D}_{\pi_\theta})$, $\theta_{i+1} = \theta_i + \alpha \nabla_{\theta} J(\theta) \big|_{\theta=\theta_i}$, $\theta_i \rightarrow \theta^*$
- We control the learning rate in parameter space, but due to the nonlinearity of neural networks, even a small change in parameter may lead to very different policies (output of the NN)—so a single bad step can collapse the training procedure. This makes it dangerous to use large step sizes with VPG, thus hurting its sample efficiency and induce risks.
- TRPO and a simpler algorithm PPO were invented to resolve this issue by defining the update constraint not on parameter space but directly on policy space.

Trust Region Policy Optimization (TRPO)

- Let π_θ denote a policy with parameters θ . We still have the similar goal with an additional constraint

$$\begin{aligned} \theta_{k+1} &= \arg \max_{\theta} L(\theta_k, \theta) \\ \text{s.t. } \bar{D}_{KL}(\theta \parallel \theta_k) &\leq \delta \end{aligned}$$

- where $L(\theta_k, \theta)$ is called the surrogate advantage, we use the old parameters to evaluate each action, and change the policy to increase the chance of “good” actions and decrease the chance of “bad” ones

$$L(\theta_k, \theta) = \mathbb{E}_{s, a \sim \pi_{\theta_k}} \left[\frac{\pi_\theta(a | s)}{\pi_{\theta_k}(a | s)} A^{\pi_{\theta_k}(s, a)} \right]$$

- $\bar{D}_{KL}(\theta \parallel \theta_k)$ is an average KL-divergence between policies across states visited by the old policy. In other words, when parameters change, the chance of performing an action should not change much

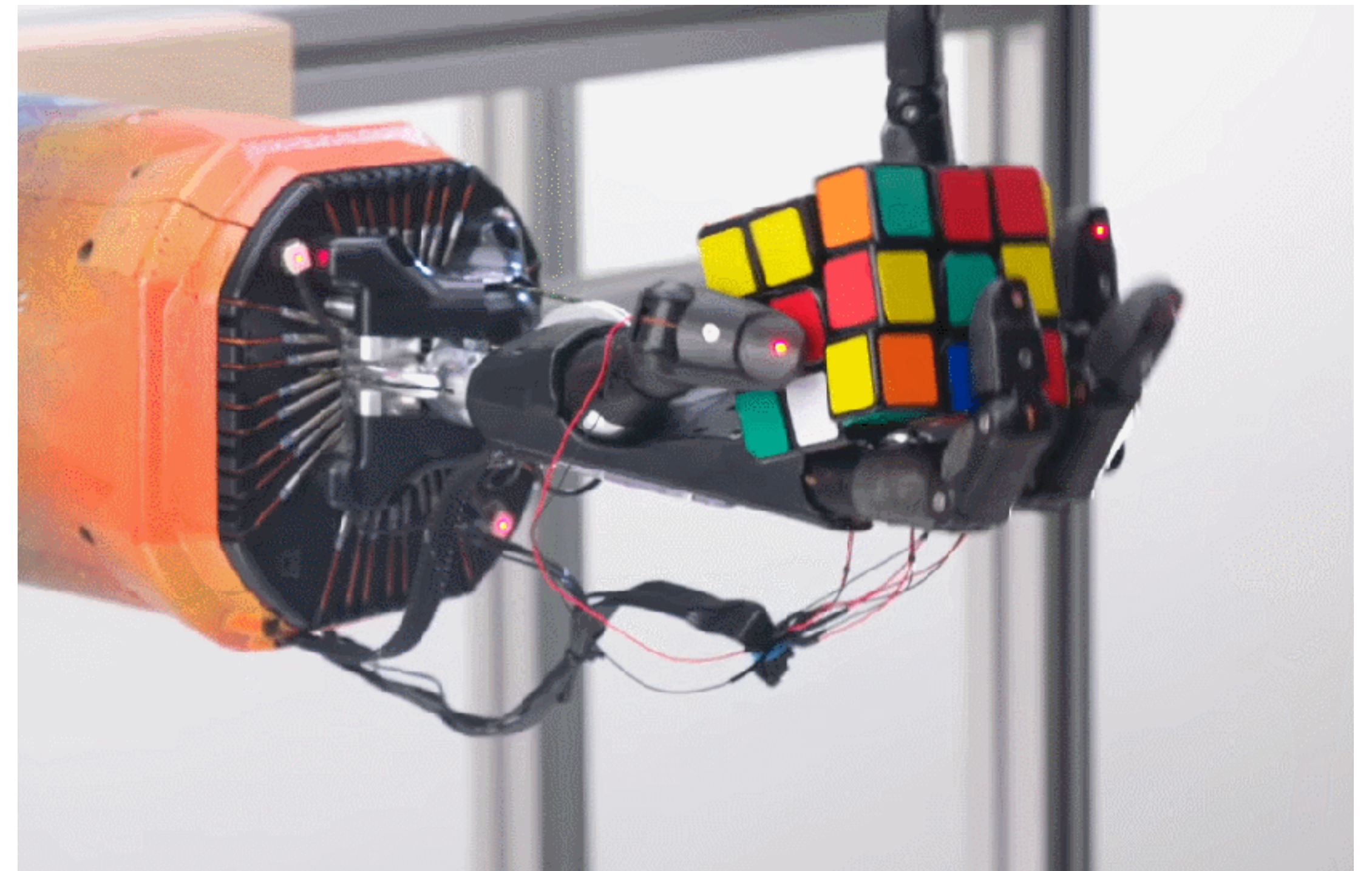
$$\bar{D}_{KL}(\theta \parallel \theta_k) = \mathbb{E}_{s \sim \pi_{\theta_k}} \left[D_{KL}(\pi_\theta(\cdot | s) \parallel \pi_{\theta_k}(\cdot | s)) \right]$$

Kullback–Leibler divergence $D_{KL}(P \parallel Q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$

Proximal Policy Optimization (PPO)

- PPO is motivated by the same question as TRPO: how can we take the biggest possible improvement step on a policy using the data we currently have, without stepping so far that we accidentally cause performance collapse?
- Where TRPO tries to solve this problem with a complex K-L divergence (second order) method, PPO is a family of first-order methods that use a few other tricks to keep new policies close to old.
- PPO methods are significantly simpler to implement, and empirically seem to perform at least as well as TRPO.
- PPO has been widely used in practice for its fast speed and simple implementation

Famous tests using PPO



PPO algorithm

- Still the same goal

$$\theta_{k+1} = \arg \max_{\theta} \mathbb{E}_{s, a \sim \pi_{\theta_k}} \left[L(s, a, \theta_k, \theta) \right]$$

- But incorporate the constraint into the objective function with a simple clip g

$$L(s, a, \theta_k, \theta) = \min \left(\frac{\pi_{\theta}(a | s)}{\pi_{\theta_k}(a | s)} A^{\pi_{\theta_k}(s, a)}, \quad g_{\epsilon} \left(A^{\pi_{\theta_k}(s, a)} \right) \right)$$

$$\text{where } g_{\epsilon}(A) = \begin{cases} (1 + \epsilon)A & A \geq 0 \quad (\text{good action}) \\ (1 - \epsilon)A & A < 0 \quad (\text{bad action}) \end{cases}$$

- For good actions, π_{θ} increase with a ceiling at $(1 + \epsilon)\pi_{\theta_k}$
- For bad actions, π_{θ} decreases with a floor at $(1 - \epsilon)\pi_{\theta_k}$

Different ways to add constraints

- VPG

- $\theta_{i+1} = \theta_i + \alpha \nabla_{\theta} J(\theta) |_{\theta=\theta_i}$
- Use a learning rate to constrain the change of θ in π_{θ}

- PPO

- $L(s, a, \theta_k, \theta) = \min \left(\frac{\pi_{\theta}(a | s)}{\pi_{\theta_k}(a | s)} A^{\pi_{\theta_k}(s, a)}, g_{\epsilon}(A^{\pi_{\theta_k}(s, a)}) \right), g_{\epsilon}(A) = \begin{cases} (1 + \epsilon)A & A \geq 0 \\ (1 - \epsilon)A & A < 0 \end{cases}$
- Bound the averaged ratio of π_{θ} and π_{θ_k} for each action

- TRPO

- $\bar{D}_{KL}(\theta || \theta_k) = \mathbb{E}_{s \sim \pi_{\theta_k}} \left[D_{KL}(\pi_{\theta}(\cdot | s) || \pi_{\theta_k}(\cdot | s)) \right]$
- Use KL divergence to bound the change of the policy distribution π_{θ} from π_{θ_k} . We may still have a large spike of change sometimes.

They are all on-policy approach.

Constrained Markov Decision Process

- Augment the MDP with a set C of auxiliary cost functions, C_1, \dots, C_m (with each one a function $C_i : S \times A \times S \rightarrow \mathbb{R}$ mapping transition tuples to costs, like the usual reward), and limits d_1, \dots, d_m . Let $J_{C_i}(\pi)$ denote the expected discounted return of policy π with respect to cost function

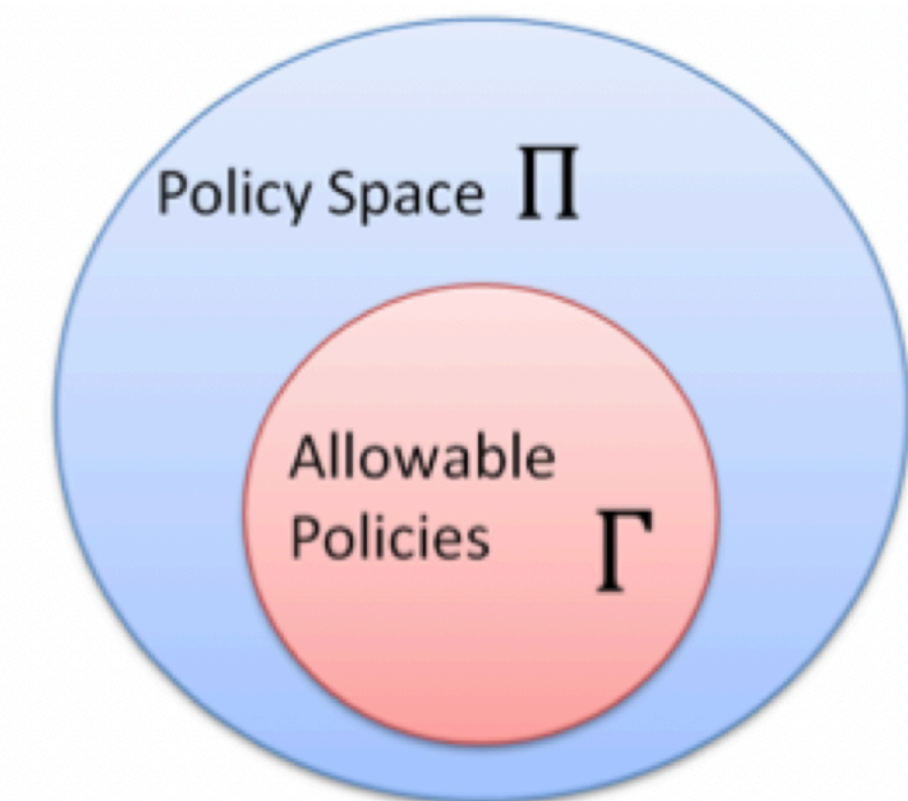
$$C_i : J_{C_i}(\pi) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t C_i(s_t, a_t, s_{t+1}) \right].$$

- The set of feasible stationary policies for a CMDP is then

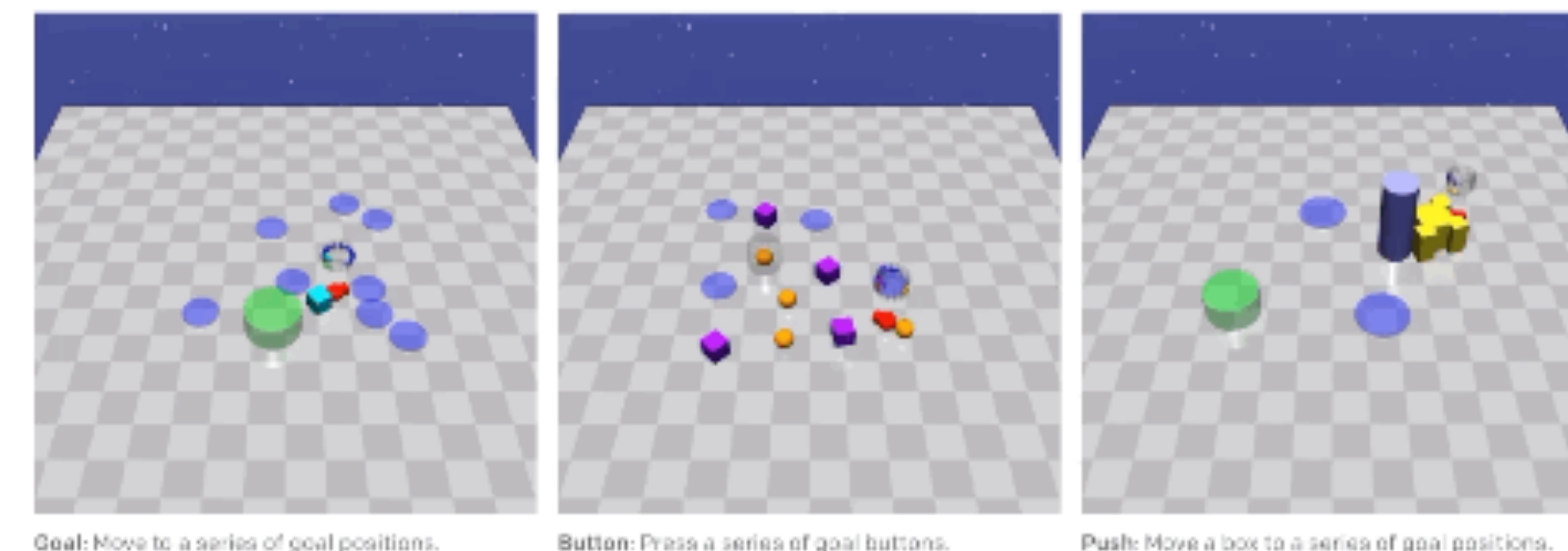
$$\Pi_C \doteq \left\{ \pi \in \Pi : \forall i, J_{C_i}(\pi) \leq d_i \right\}$$

- and the reinforcement learning problem in a CMDP is

$$\pi^* = \arg \max_{\pi \in \Pi_C} J(\pi)$$



Point is a simple robot constrained to the 2D plane, with one actuator for turning and another for moving forward or backward. Point has a front-facing small square which helps with the Push task.



Primal-dual method 1.0

- Recall the constrained objective function in CMDP with one constraint:

$$\pi^* = \arg \max_{\pi} J(\pi) \quad \text{maximizing the reward return}$$

$$s.t. \quad J_C(\pi) \leq d \quad \text{satisfying the constraint}$$

- It has an equivalent Lagrangian formulation:

$$\pi^*, \lambda^* = \arg \min_{\lambda \geq 0} \max_{\pi} L(\pi, \lambda)$$

$$L(\pi, \lambda) = J(\pi) - \lambda(J_C(\pi) - d)$$

- Dual problem: the outer minimization over the dual variable $\lambda \in \mathbb{R}_{\geq 0}$

Primal-dual method 1.0

- The Lagrangian formulation: $\pi^*, \lambda^* = \arg \min_{\lambda \geq 0} \max_{\pi} L(\pi, \lambda)$

$$L(\pi, \lambda) = J(\pi) - \lambda(J_C(\pi) - d)$$

- Primal-dual (Lagrangian) algorithm:

1. Fix λ , optimize π : $\pi_{i+1} = \arg \max_{\pi} L(\pi, \lambda_i)$

2. Fix π , optimize λ : $\lambda_{i+1} = \arg \min_{\lambda \geq 0} L(\pi_{i+1}, \lambda)$

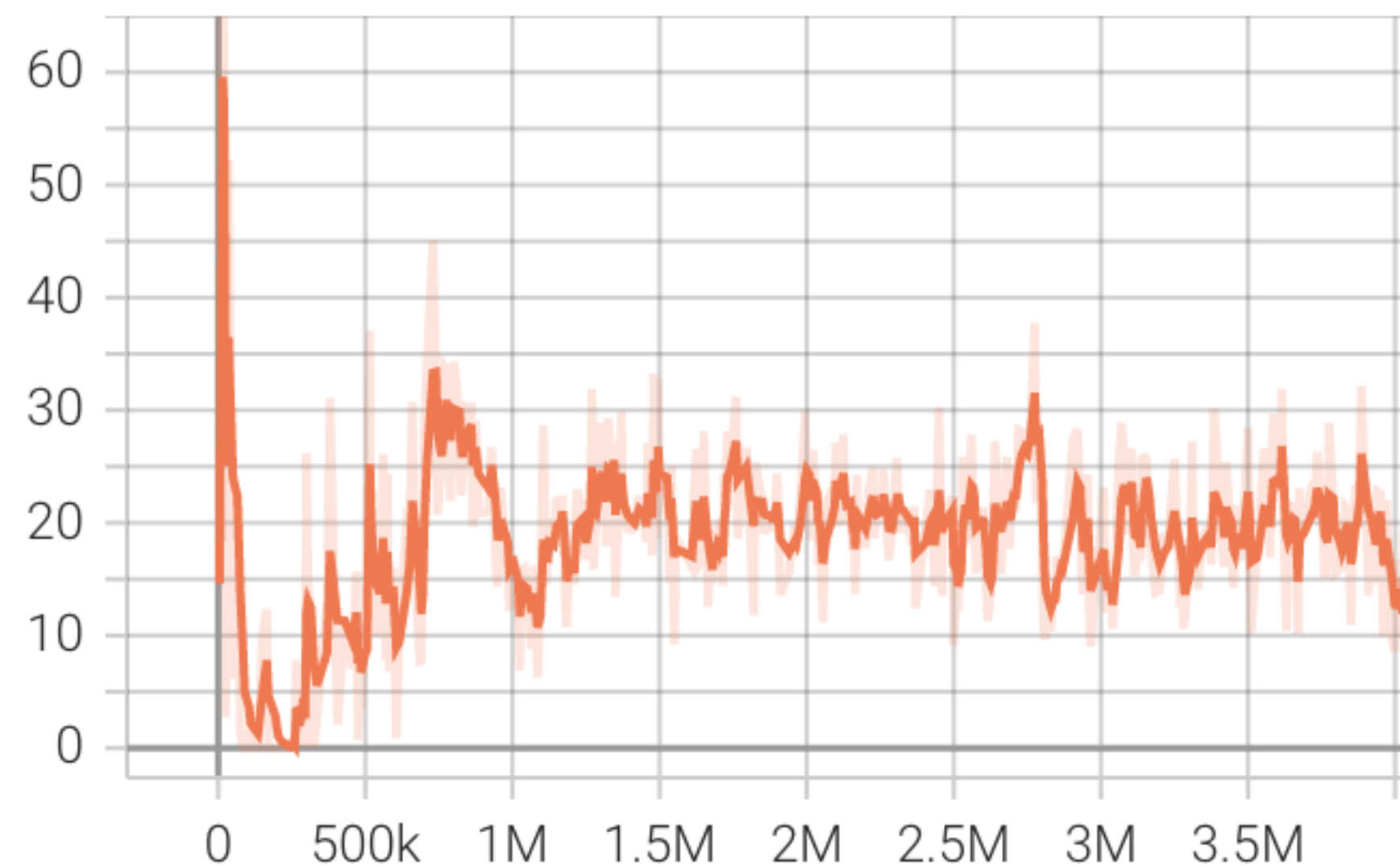
- Phase 1) can be done via any policy gradient-based RL algorithms
- Phase 2) is usually done via one-step gradient descend

$$\lambda_{i+1} = \max \left(0, \lambda_i - \eta \nabla_{\lambda} L(\pi_{i+1}, \lambda_i) \right) = \max \left(0, \lambda_i + \eta (J_C(\pi_{i+1}) - d) \right)$$

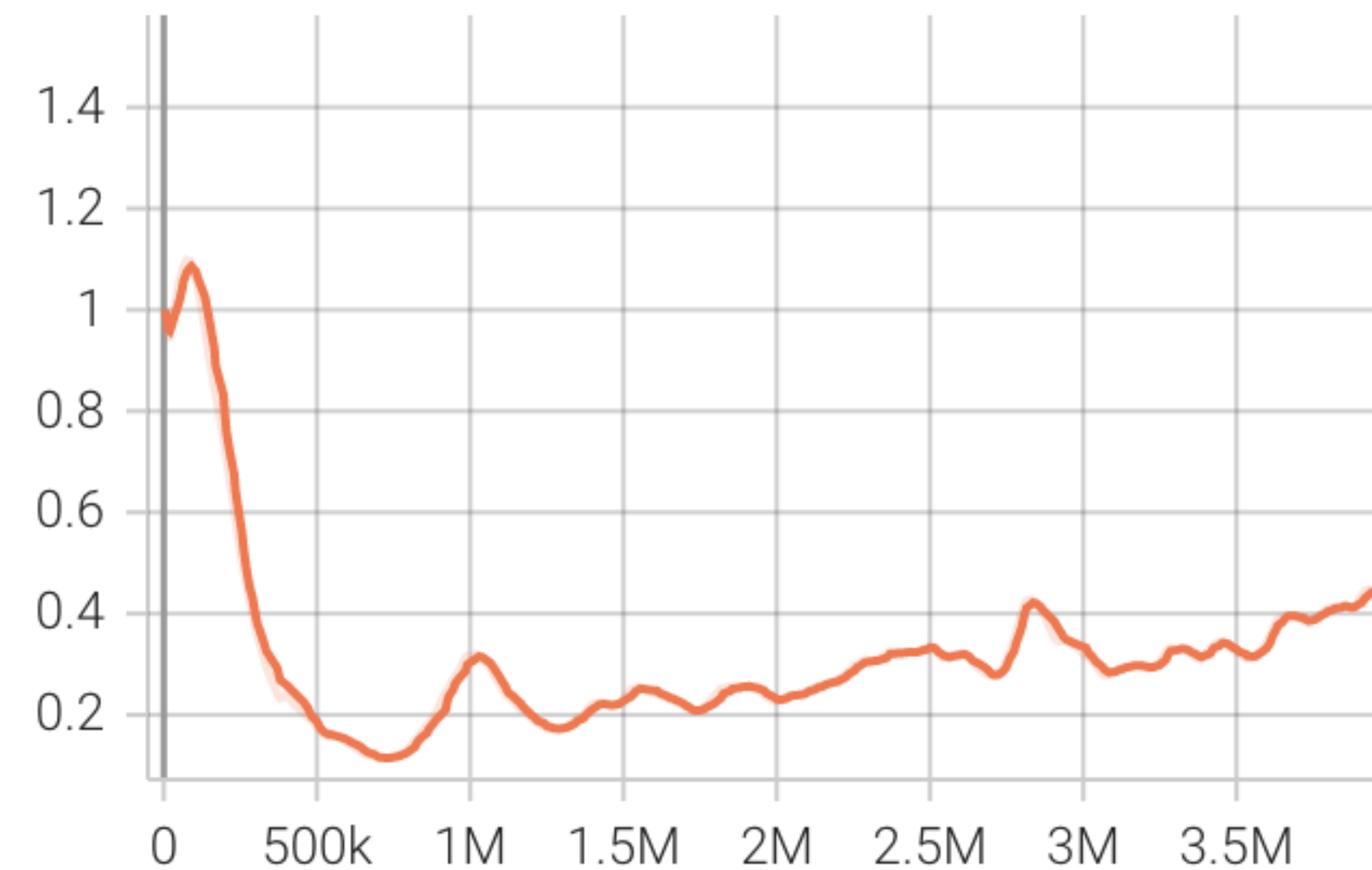
Primal-dual method 1.0

- Example: using PPO to optimize $\pi_{i+1} = \arg \max_{\pi} J(\pi) - \lambda(J_C(\pi) - d)$
- PPO-Lagrangian algorithm

EpCost $J_C(\pi)$
tag: EpCost



Penalty λ
tag: Penalty



Primal-dual method 1.0

- What's the problem?
 - The dual problem is only approximately solved: $\lambda \rightarrow 0$ when $J_C(\pi)$ satisfies the constraint, $\lambda \rightarrow +\infty$ we fail

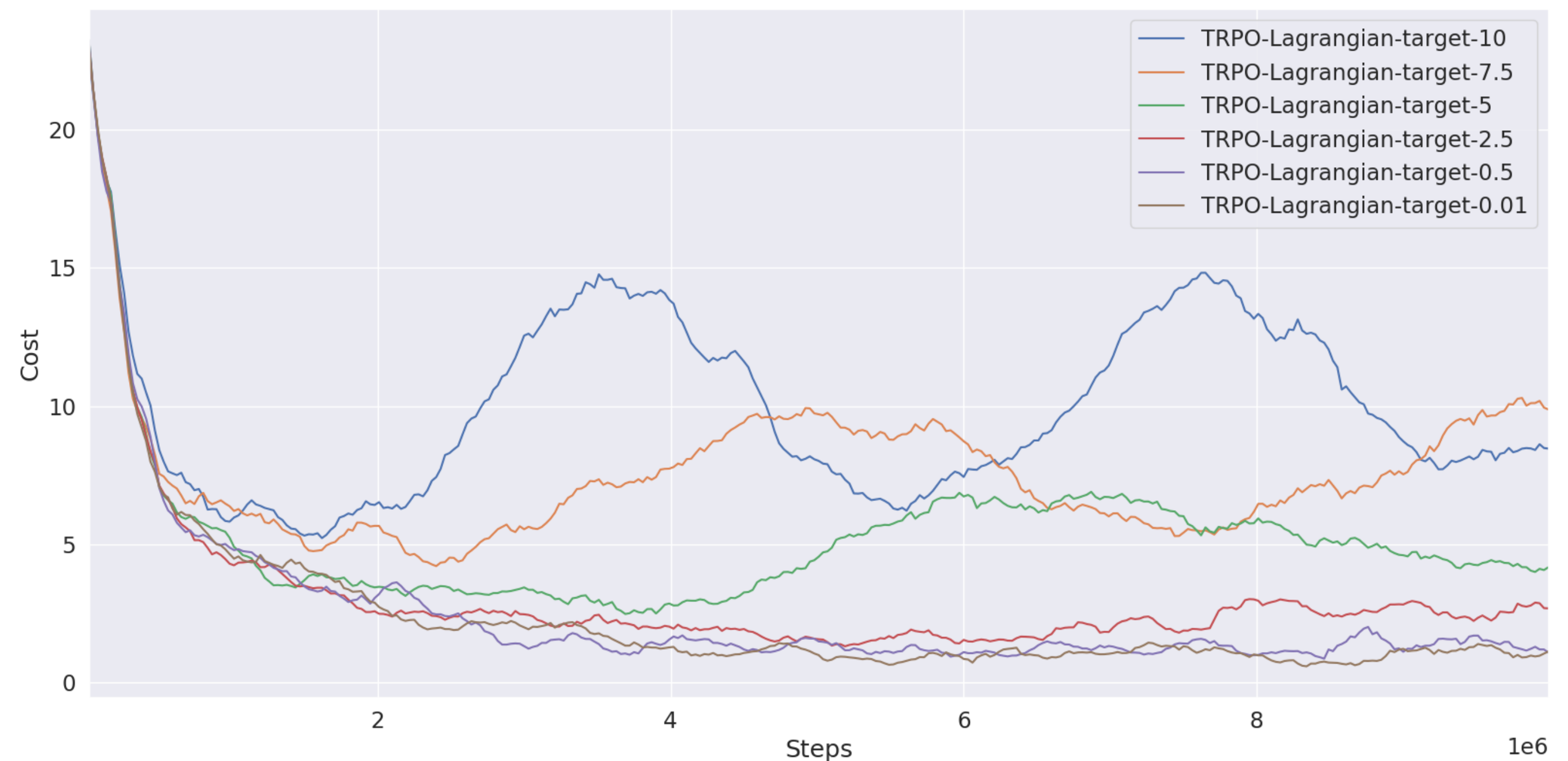
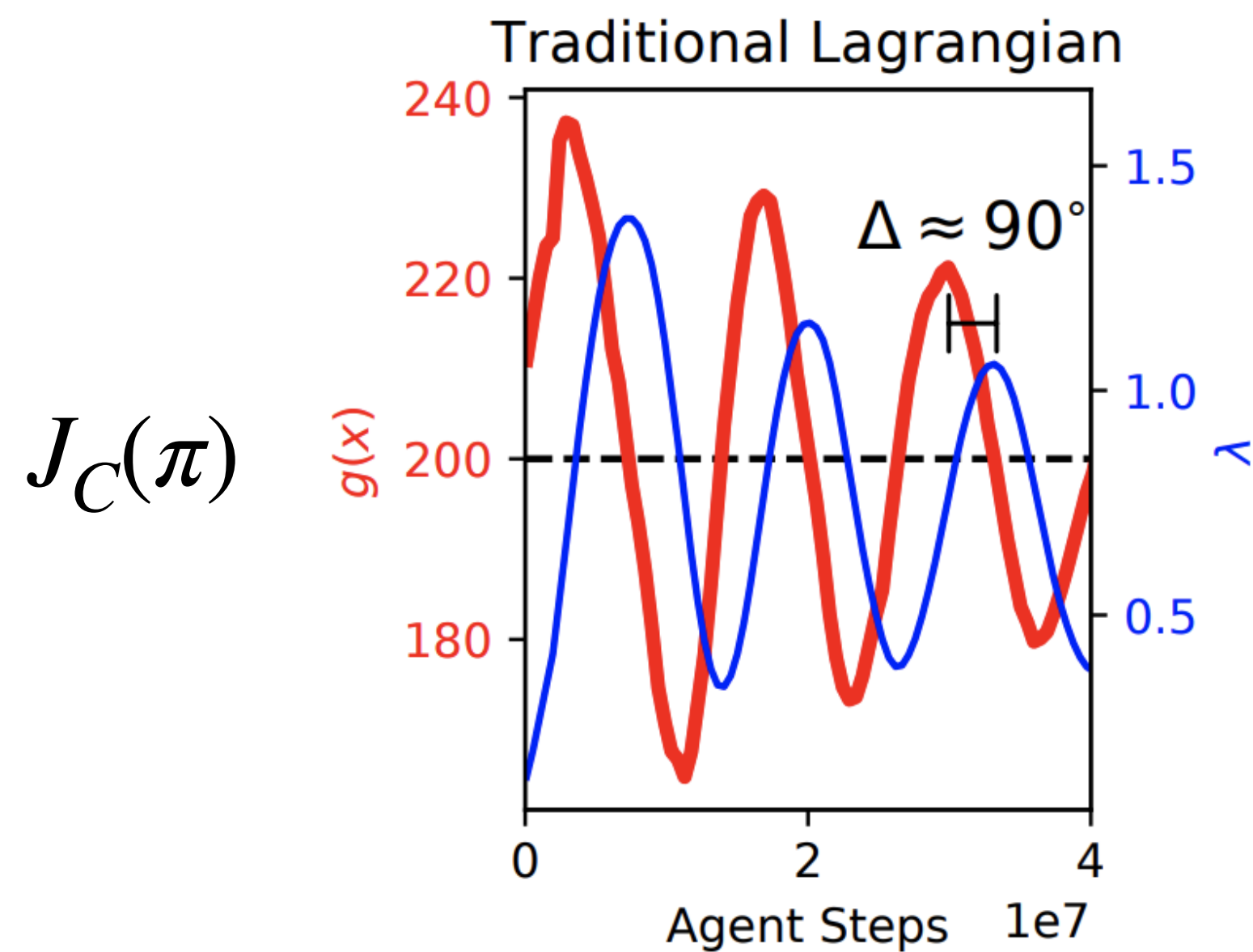
$$\pi^*, \lambda^* = \arg \min_{\lambda \geq 0} \max_{\pi} J(\pi) - \lambda(J_C(\pi) - d)$$

- We need to select a proper learning rate for λ

$$\lambda_{i+1} = \max \left(0, \lambda_i + \eta(J_C(\pi_{i+1}) - d) \right)$$

Primal-dual method 2.0

- What's the problem?
 - Improper learning rate can lead to oscillation training behavior and phase shift between λ and $J_C(\pi)$



Primal-dual method 2.0

- How to solve the unstable training problem?
- Denote $e_i = J_C(\pi_{i+1}) - d$, we could observe that

$$\lambda_{i+1} = \lambda_i + \eta e_i = \lambda_{i-1} + \eta(e_i + e_{i-1}) = \lambda_1 + \eta \sum_{j=1}^i e_j$$

- The dual problem updating is an integral controller for λ !
- Then the oscillation, phase shift problems could be explained.
- We can regard the dual variable updating as a control system and design a PID controller to solve λ .

Primal-dual method 2.0

- PID-Lagrangian algorithm

Algorithm 2 PID-Controlled Lagrange Multiplier

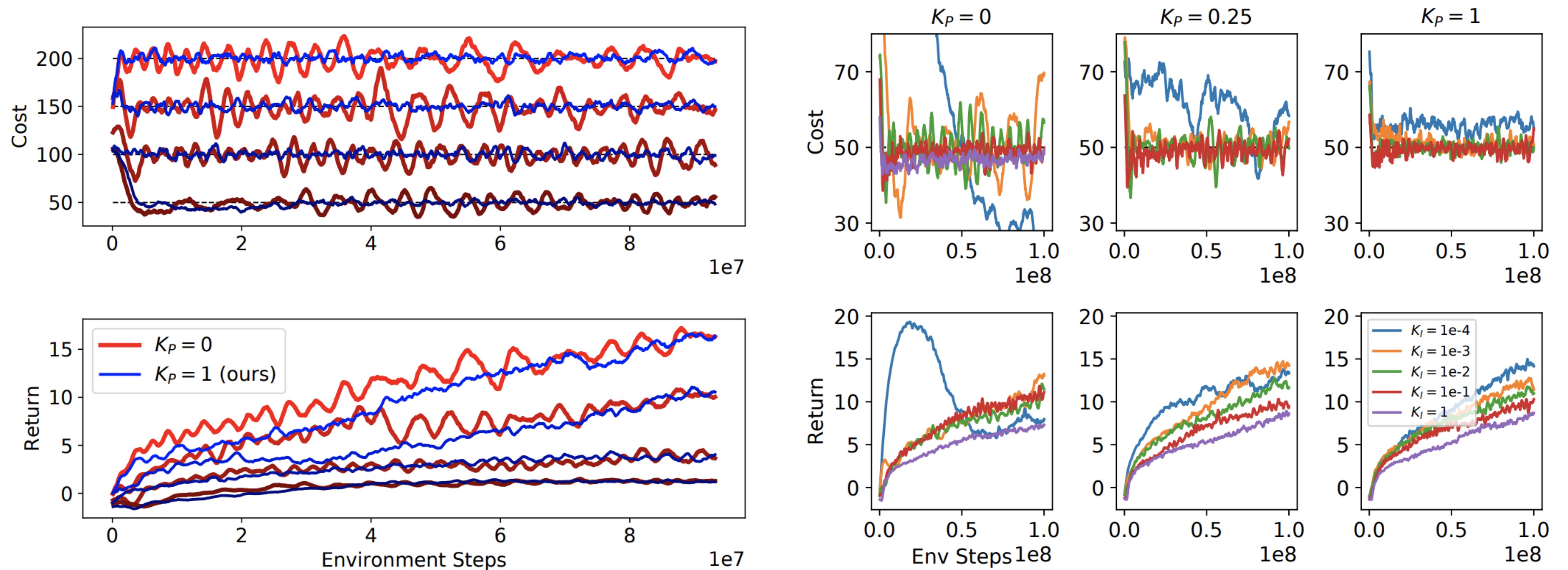
```
1: Choose tuning parameters:  $K_P, K_I, K_D \geq 0$ 
2: Integral:  $I \leftarrow 0$ 
3: Previous Cost:  $J_{C,prev} \leftarrow 0$ 
4: repeat at each iteration  $k$ 
5:   Receive cost  $J_C$ 
6:    $\Delta \leftarrow J_C - d$ 
7:    $\partial \leftarrow (J_C - J_{C,prev})_+$ 
8:    $I \leftarrow (I + \Delta)_+$ 
9:    $\lambda \leftarrow (K_P\Delta + K_I I + K_D\partial)_+$ 
10:   $J_{C,prev} \leftarrow J_C$ 
11:  return  $\lambda$ 
```

Algorithm 1**Constraint-Controlled Reinforcement Learning**

```
1: procedure CONSTRAINED RL( $\pi_{\theta_0}(\cdot|s), d$ )
2:   Initialize control rule (as needed)
3:    $\mathcal{J}_C \leftarrow \{\}$   $\triangleright$  cost measurement history
4:   repeat
5:     Sample environment:  $\triangleright$  a minibatch
6:      $a \sim \pi(\cdot|s; \theta), s' \sim T(s, a),$ 
7:      $r \sim R(s, a, s'), c \sim C(s, a, s')$ 
8:     Apply feedback control:
9:     Store sample estimate  $\hat{J}_C$  into  $\mathcal{J}_C$ 
10:     $\lambda \leftarrow h(\mathcal{J}_C, d), \lambda \geq 0$ 
11:    Update  $\pi$  by RL:  $\triangleright$  by Lagrangian objective
12:    Update critics,  $V_\phi(s), V_{C,\psi}(s)$   $\triangleright$  if using
13:     $\nabla_\theta \mathcal{L} = \frac{1}{1+\lambda} \left( \nabla_\theta \hat{J}(\pi_\theta) - \lambda \nabla_\theta \hat{J}_C(\pi_\theta) \right)$ 
14:  until converged
15:  return  $\pi_\theta$ 
16: end procedure
```

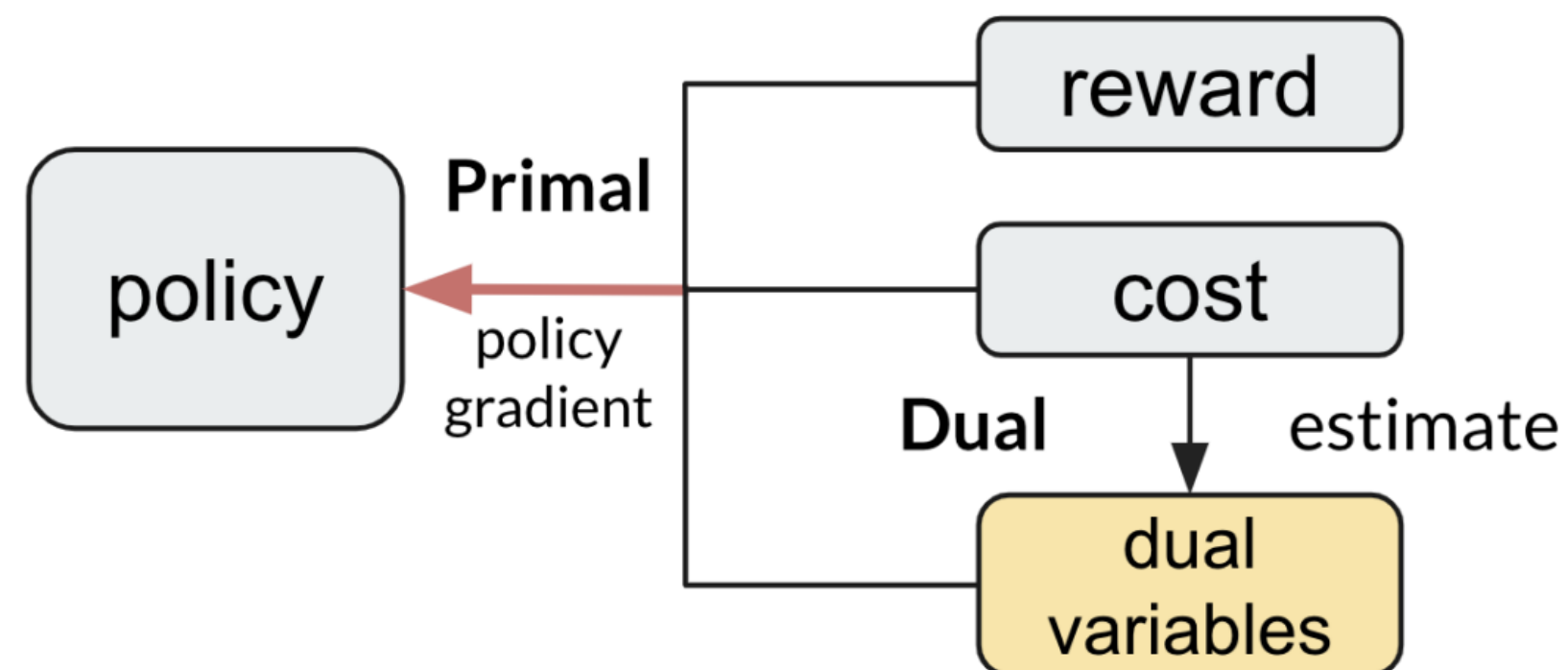
Primal-dual method 2.0

- PID-Lagrangian algorithm has better training stability



Primal-dual method 3.0

- The problems for primal-dual methods:
 - Lack of optimality guarantees both during training and after training
 - The policy only converges to a saddle point asymptotically
 - The primal problem is usually difficult to optimize



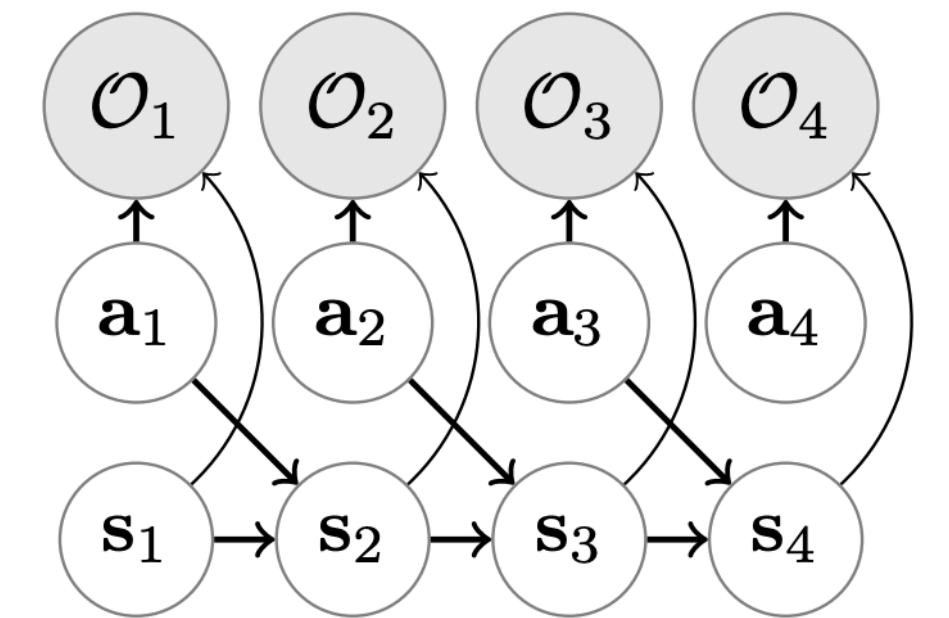
- How to overcome the theoretical drawbacks of primal-dual methods?

Variational inference approach

- Consider the safe RL problem from the probabilistic inference perspective

- Denote \mathcal{O} as the optimality variable

- $p(\mathcal{O}_t = 1 | s_t, a_t) = K_{\text{norm}} \exp(r(s_t, a_t)) \propto \exp(r(s_t, a_t))$



graphical model with optimality variables

- We have the following Evidence Lower Bound (ELBO) for likelihood:

$$\log p_{\pi}(\mathcal{O} = 1) = \log \int p(\mathcal{O} = 1 | \tau) p_{\pi}(\tau) d\tau \geq \mathbb{E}_{\tau \sim q} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] - \alpha KL(q || p_{\pi}) = J(q, \pi)$$

ELBO Derivation

- Denote the trajectory probability of a policy as:

$$p_{\pi_{\theta}}(\tau) = p(s_0) \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t) p(\theta) \quad q(\tau) = p(s_0) \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) q(a_t | s_t)$$

- q is the variational distribution. We have:

$$\log p_{\pi}(O = 1) = \log \int p(O = 1 | \tau) p_{\pi}(\tau) d\tau$$

$$= \log \mathbb{E}_{\tau \sim q} \left[\frac{p(O = 1 | \tau) p_{\pi}(\tau)}{q(\tau)} \right]$$

$$\geq \mathbb{E}_{\tau \sim q} \log \frac{p(O = 1 | \tau) p_{\pi}(\tau)}{q(\tau)} \quad \Rightarrow \text{Jensen's inequality}$$

$$= \mathbb{E}_{\tau \sim q} \log p(O = 1 | \tau) + \mathbb{E}_{\tau \sim q} \log \frac{p_{\pi}(\tau)}{q(\tau)}$$

$$\propto \underbrace{\mathbb{E}_{\tau \sim q} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]}_{\text{trajectory-wise ELBO}} - \alpha D_{\text{KL}}(q(\tau) || p_{\pi}(\tau)) = \mathcal{J}(q, \pi) = \mathbb{E}_{\tau \sim q} \left[\sum_{t=0}^{\infty} (\gamma^t r_t - \alpha D_{\text{KL}}(q(\cdot | s_t) || \pi_{\theta}(\cdot | s_t))) \right] \text{ state-wise ELBO} + \log p(\theta), \quad \forall q(a | s_t) \in \Pi_Q^{\epsilon_1}.$$

Solving Variational inference (details in the paper)

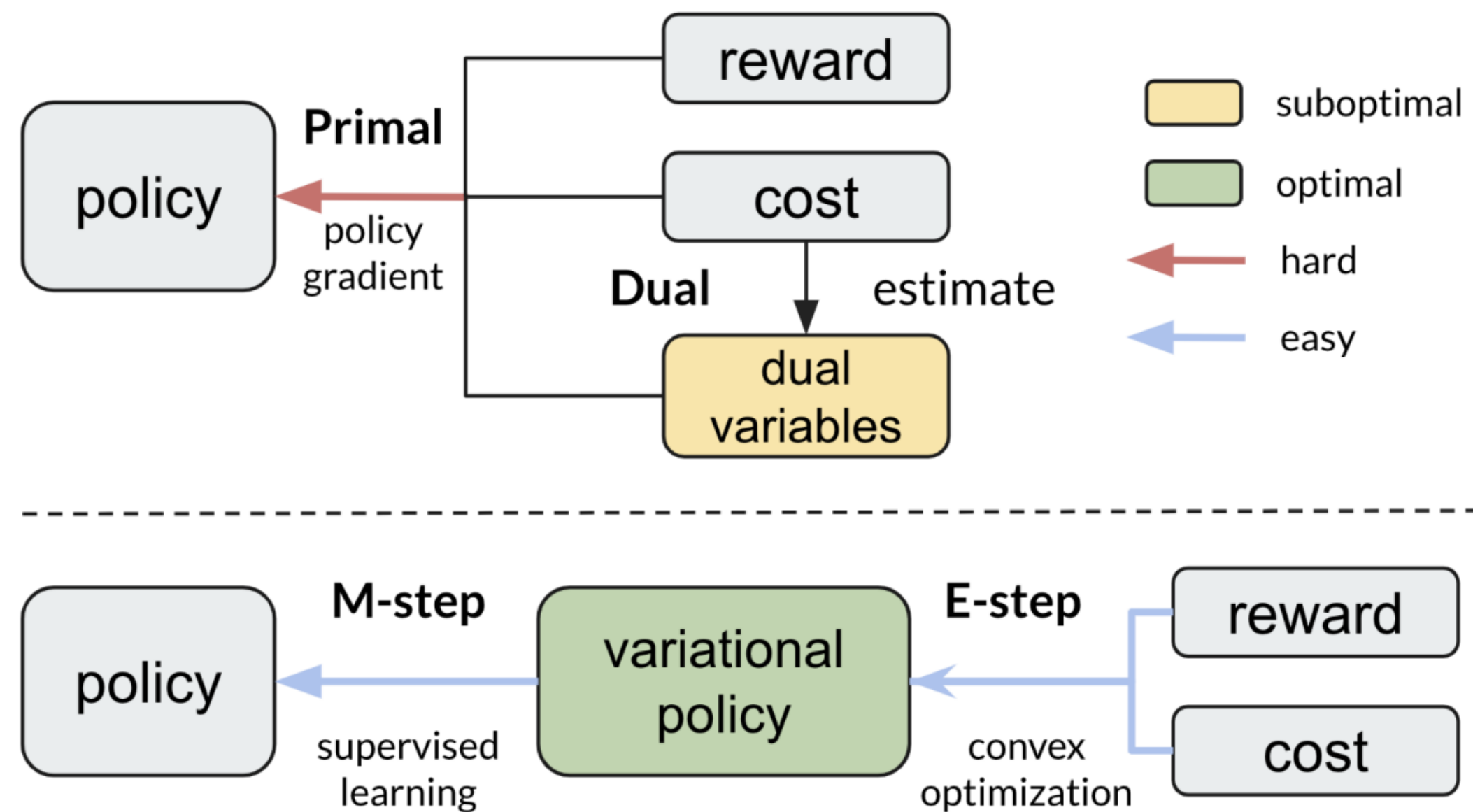
- We could improve the ELBO by Expectation-Maximization (EM) algorithm

$$J(q, \pi) = \mathbb{E}_{\tau \sim q} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] - \alpha KL(q \| p_{\pi})$$

- E-step: improve the variational distribution q
 - q can be solved analytically with optimality guarantee
- M-step: improve the policy π by minimizing the divergence $KL(q \| p_{\pi})$
 - By restricting the updated policy within a trust-region, we could achieve robustness guarantee and worst-case performance bound

Advantages

- Comparison to primal-dual method



- Robustness analysis

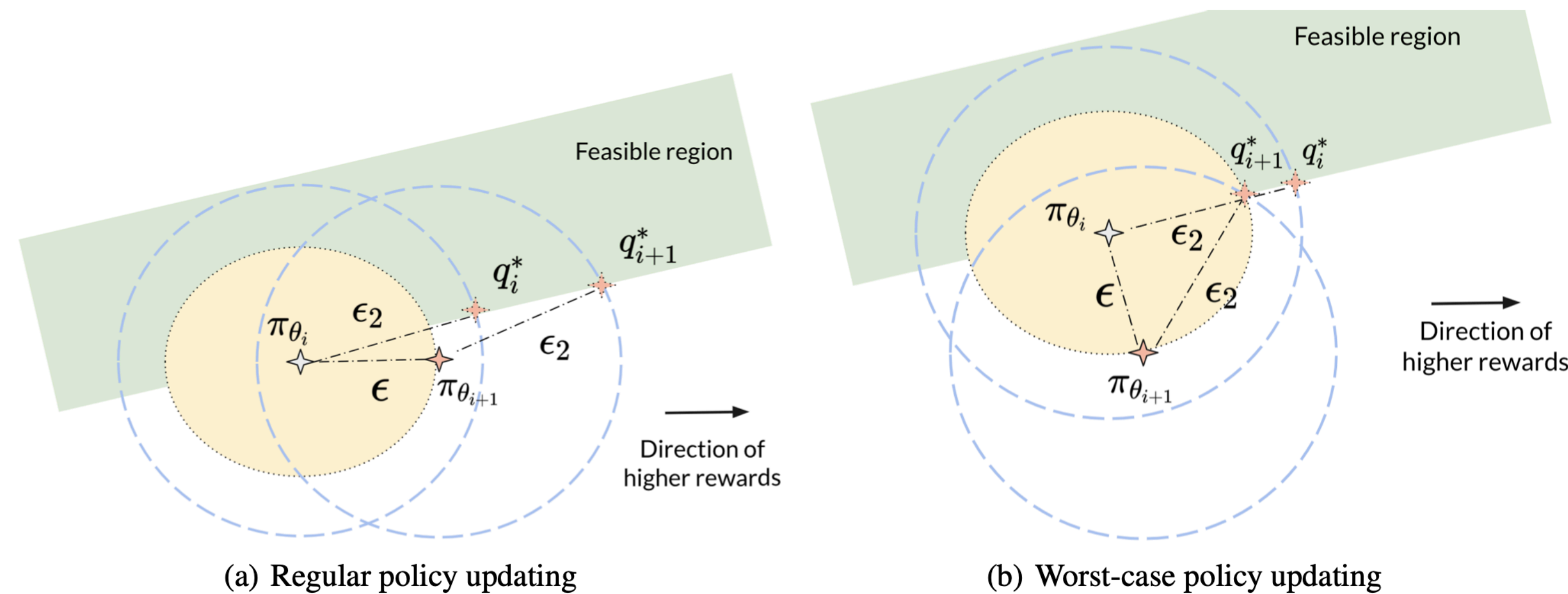


Figure 7. Illustration of the policy updating at the i -th iteration under the M-step approximation error.

- Additional benefits include: sample-efficiency (off-policy), stable training performance, optimality guarantee for **each** policy update, etc...

Worthy Reading

- Ray, Alex, Joshua Achiam, and Dario Amodei. "Benchmarking safe exploration in deep reinforcement learning." arXiv preprint arXiv:1910.01708 7 (2019).
- https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html
- Garcia, Javier, and Fernando Fernández. "A comprehensive survey on safe reinforcement learning." Journal of Machine Learning Research 16.1 (2015): 1437-1480.