

Solving Parameter-Robust Avoid Problems with Unknown Feasibility using Reinforcement Learning

Oswin So^{*†}, **Eric Yang Yu^{*†}**, Songyuan Zhang[†], Matthew Cleaveland[‡], Mitchell Black[‡], and Chuchu Fan[†]

[†]Department of Aeronautics and Astronautics, MIT

[‡]MIT Lincoln Laboratory

^{*}Equal Contribution



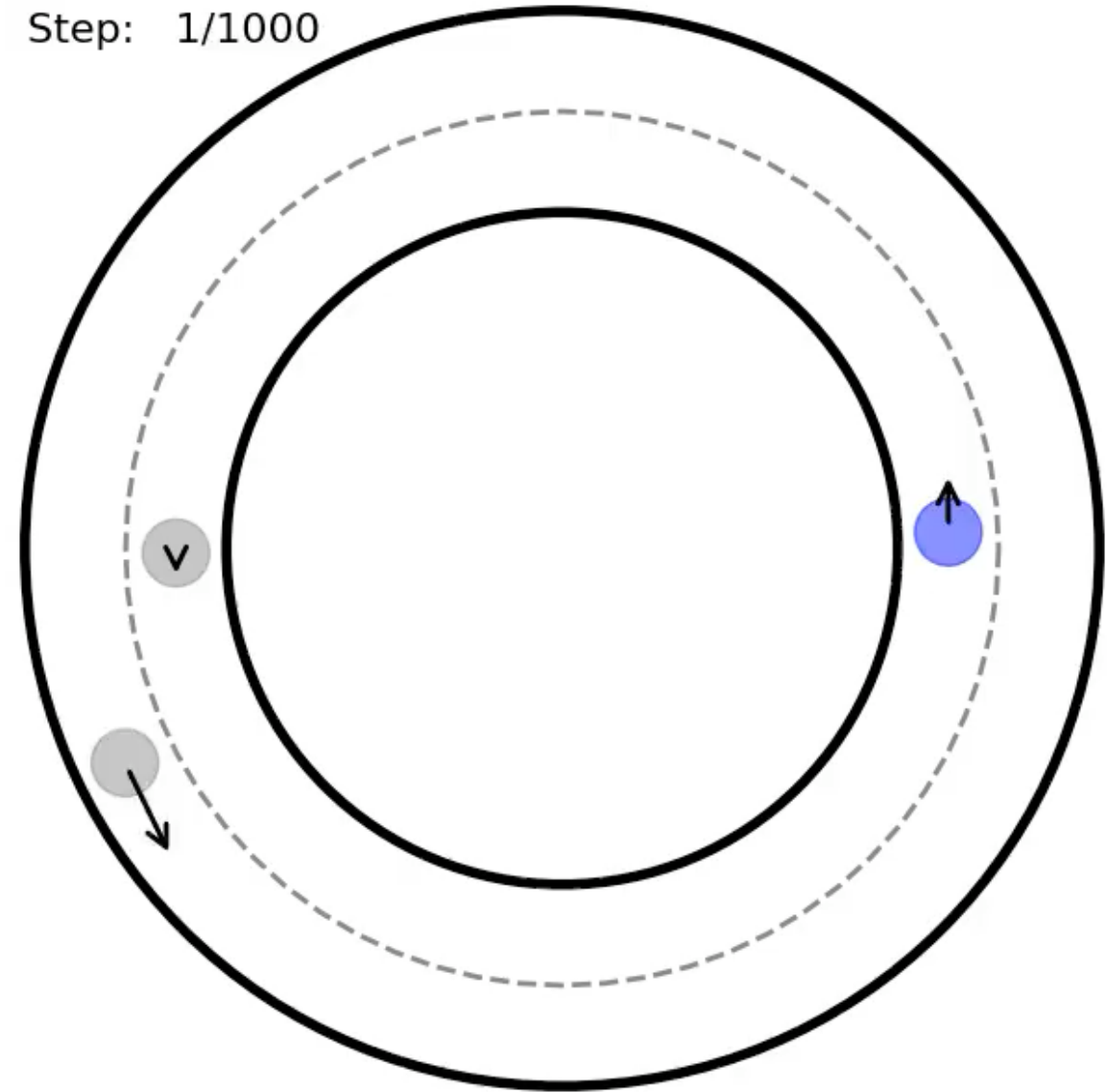
ICLR

Step: 1/1000

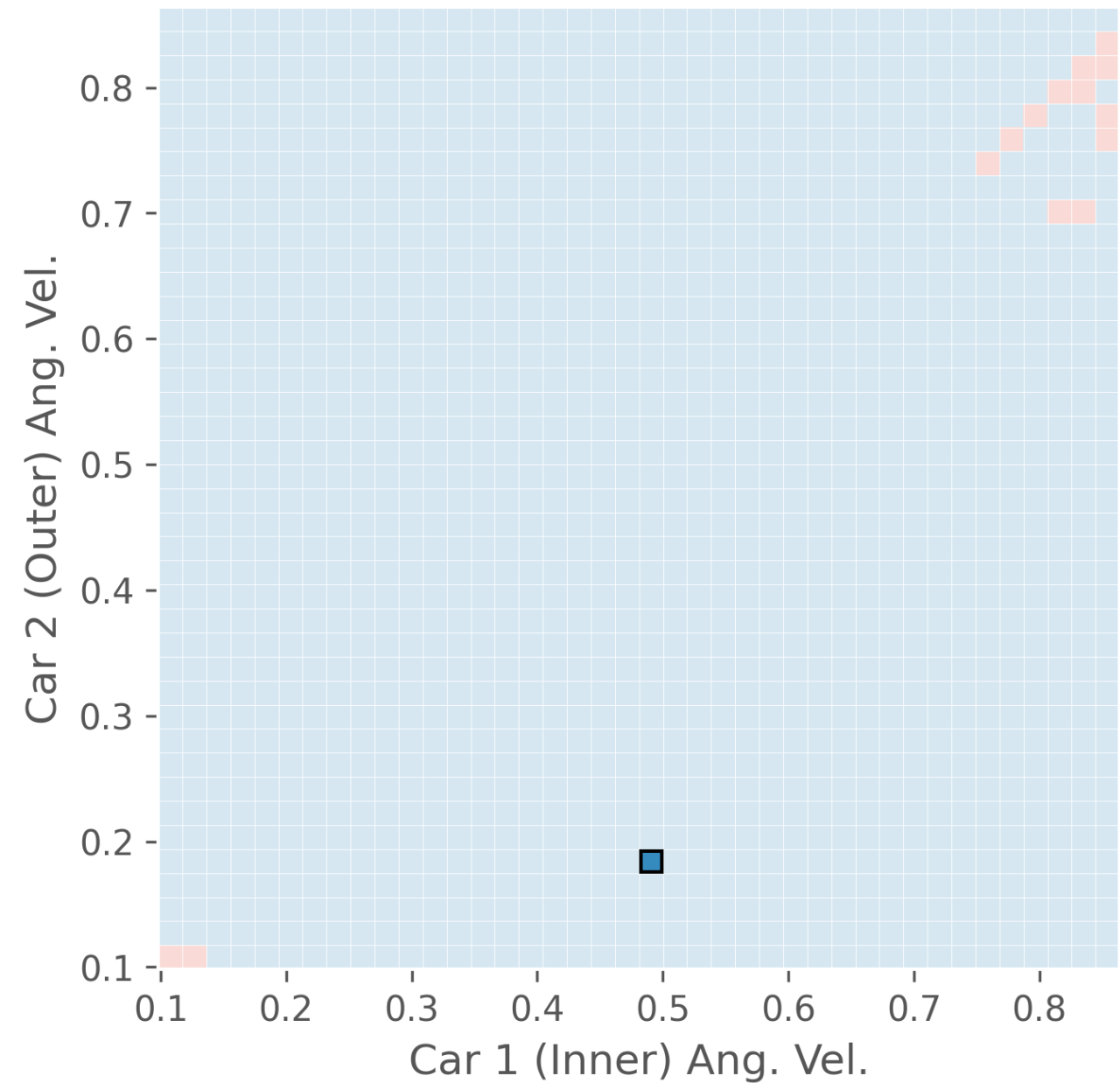
Safety Specification:

- Stay inside road
- Don't collide
- Stay within velocity limits

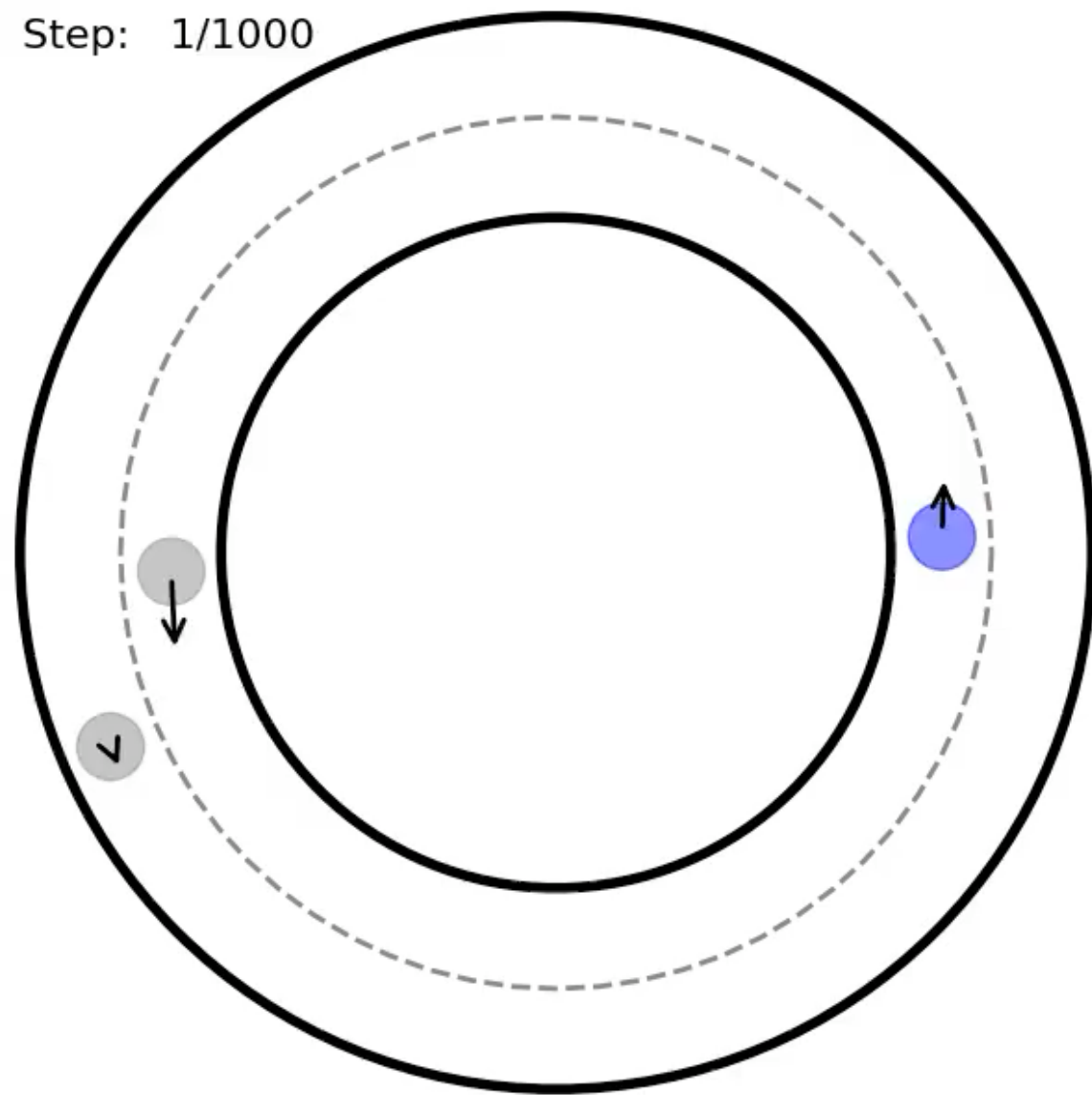
Be as robust as possible to the
other cars' velocities
(**parameters**)



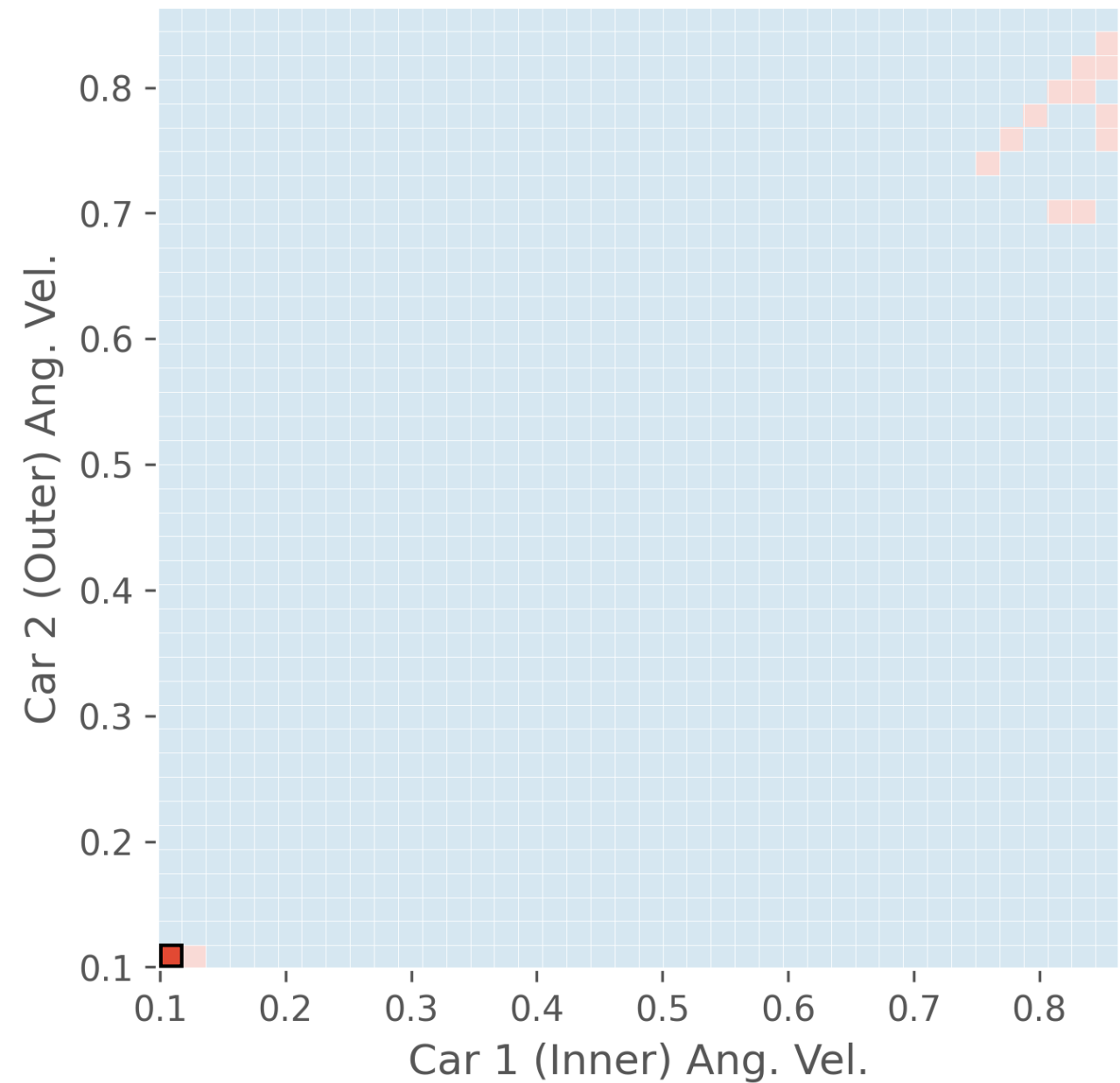
Unsafe Safe



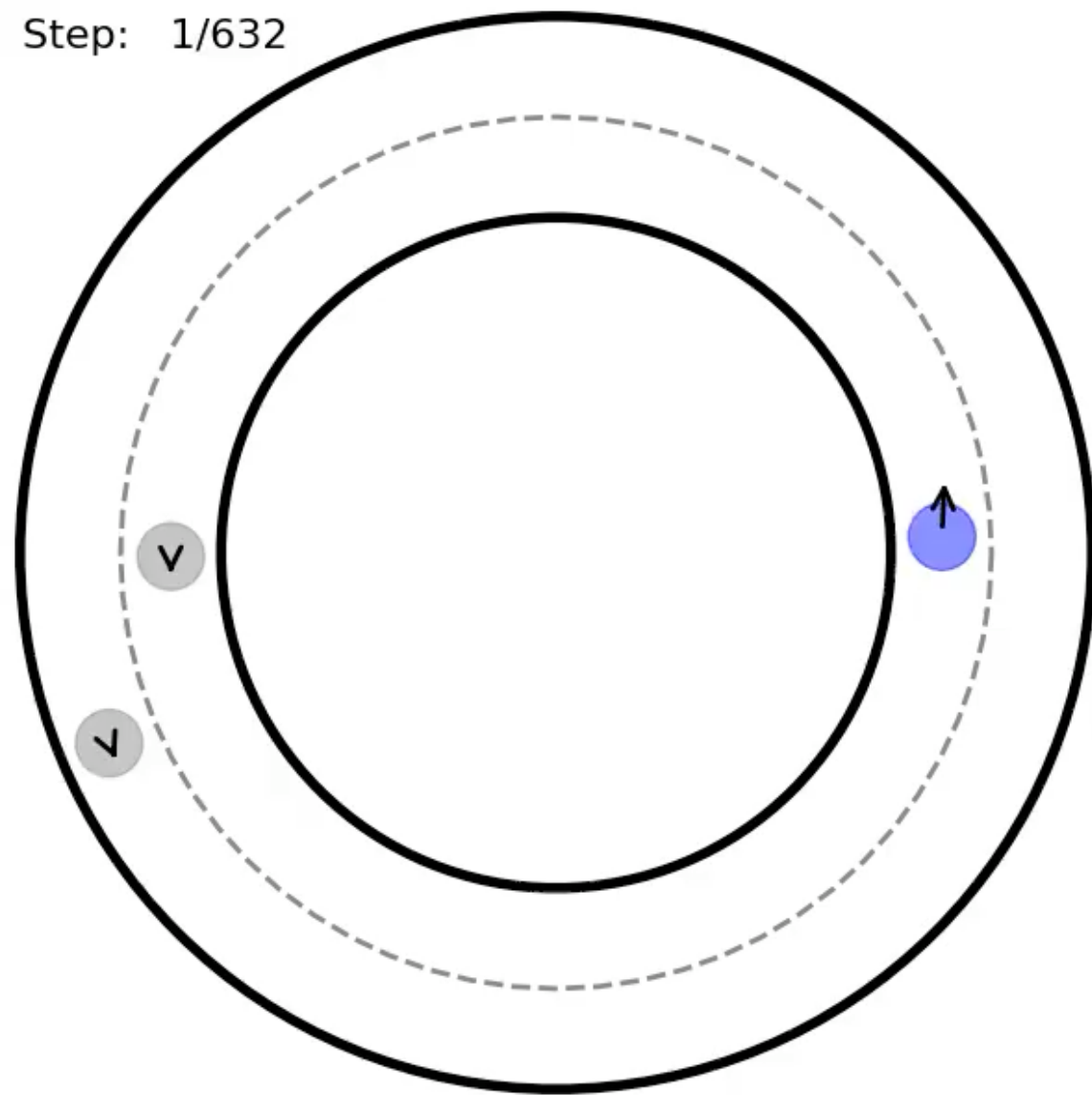
Step: 1/1000



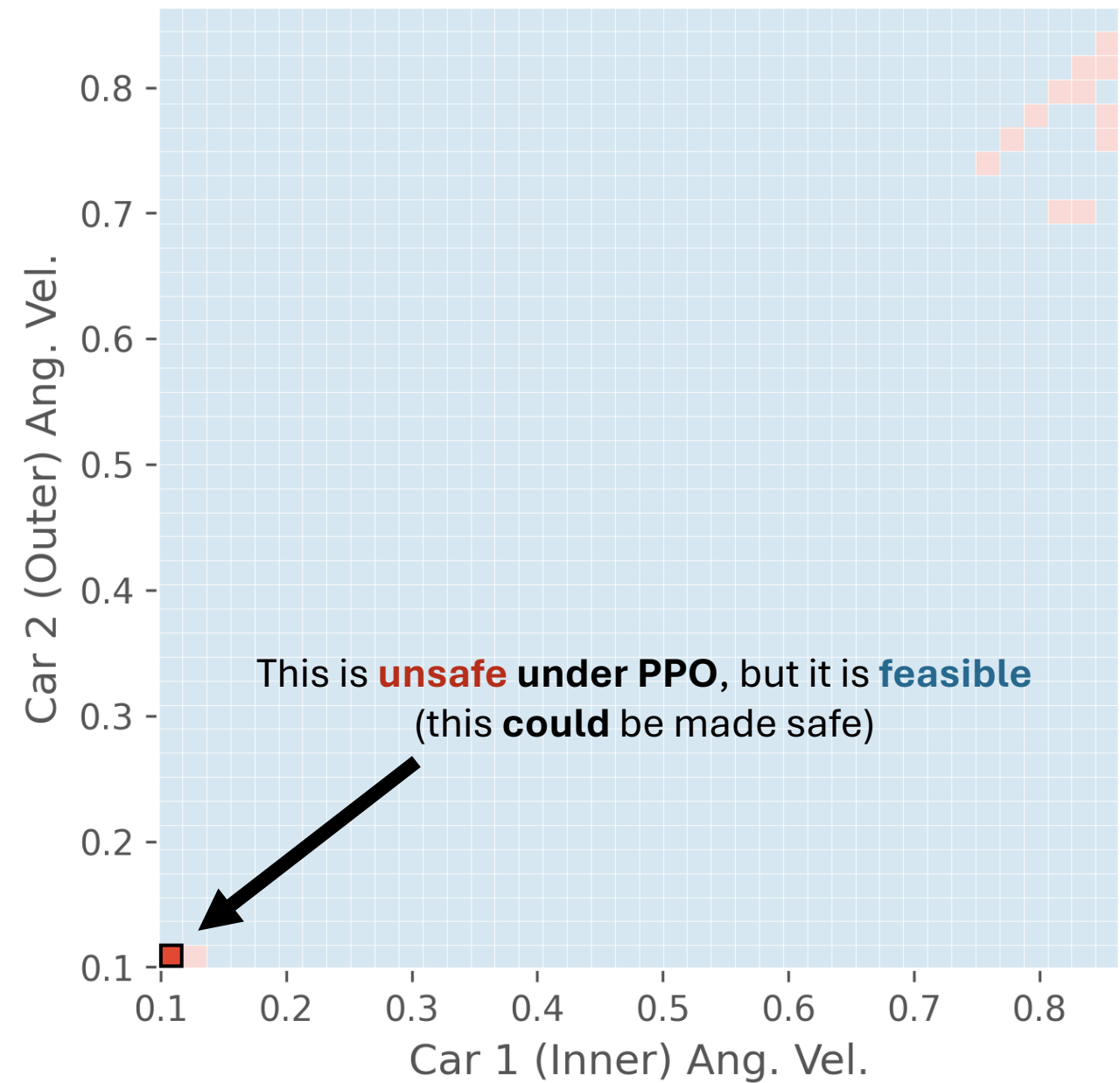
Unsafe Safe



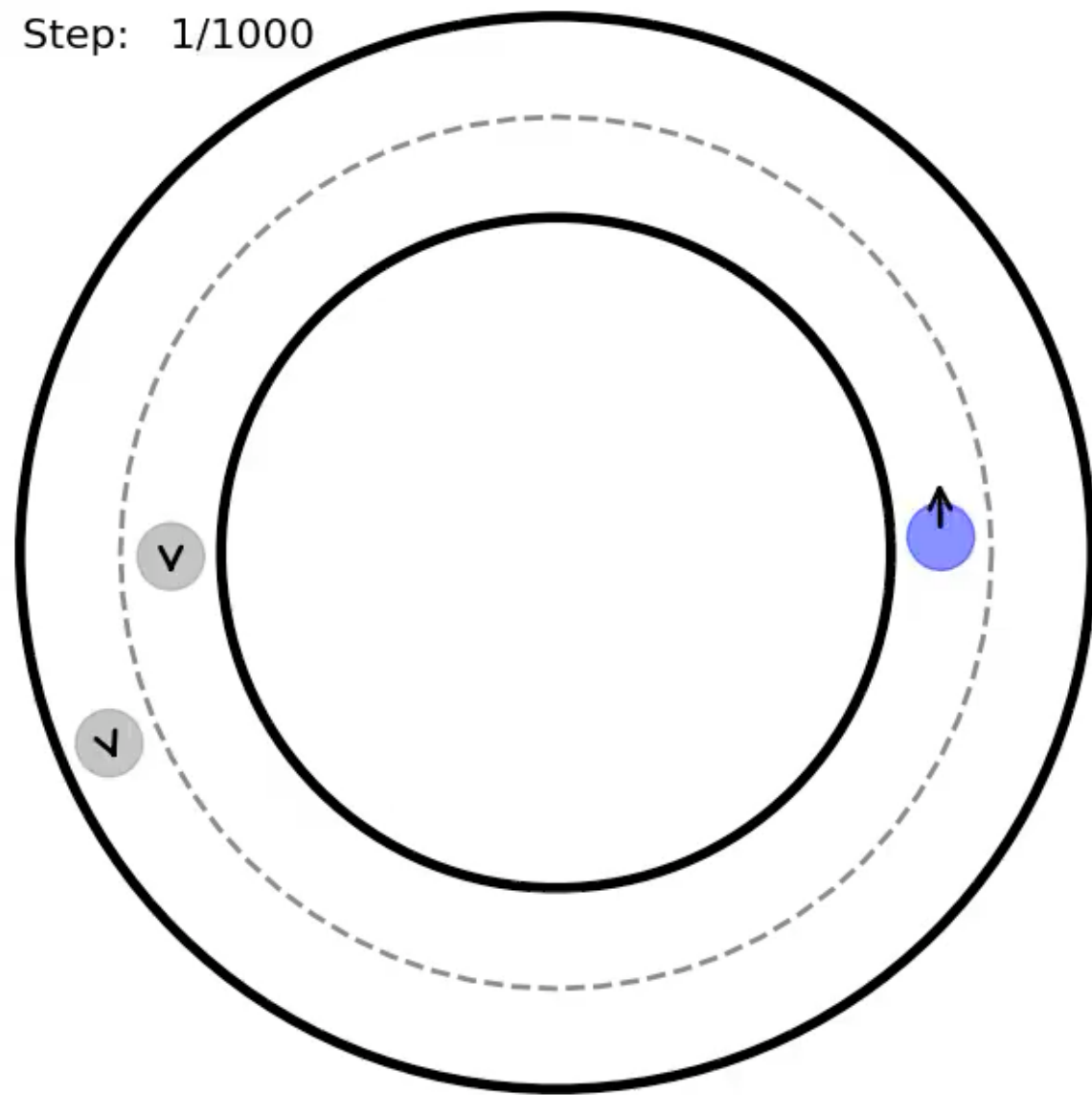
Step: 1/632



Unsafe Safe



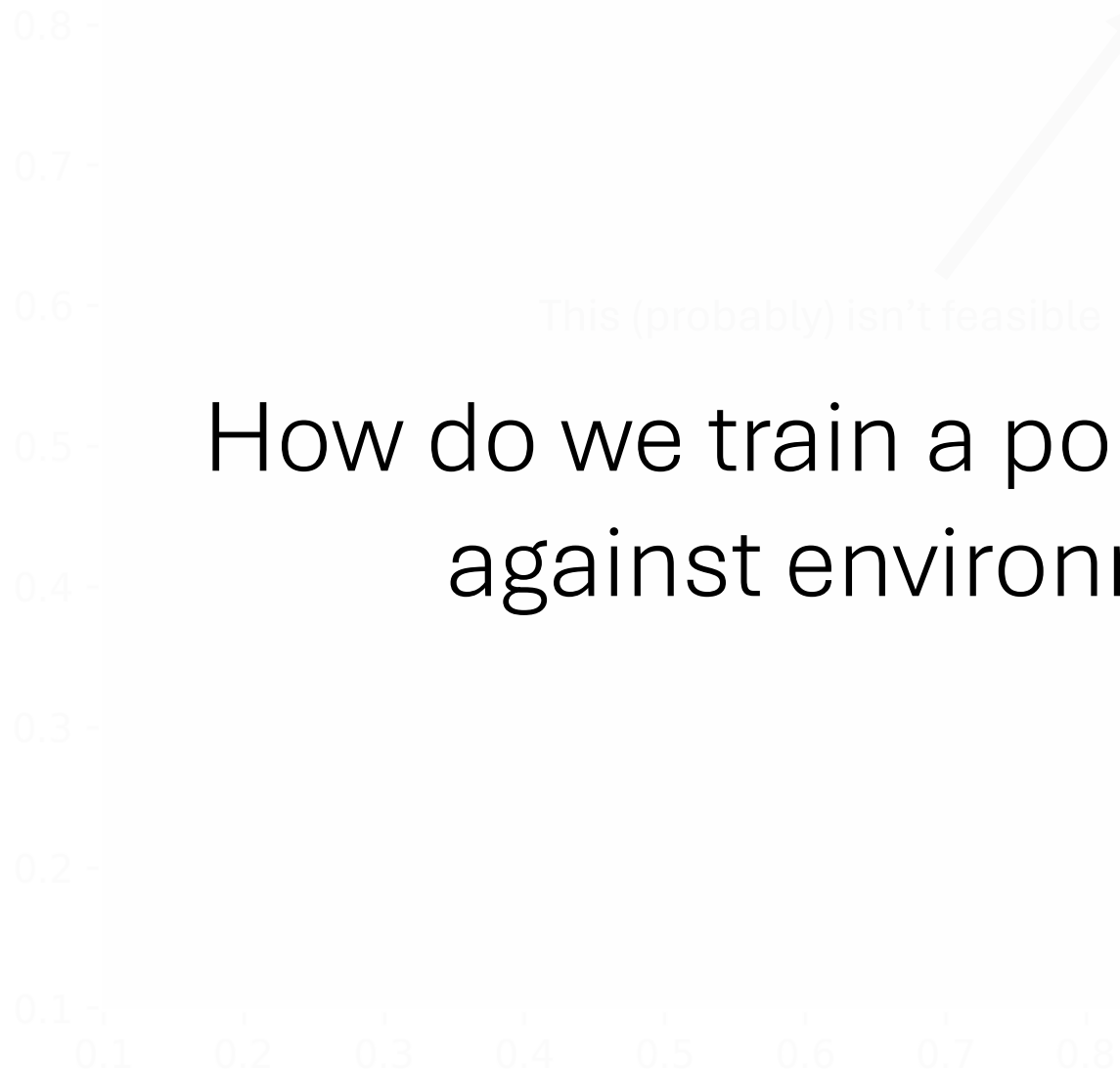
Step: 1/1000



Unsafe Safe

Step: 1/195

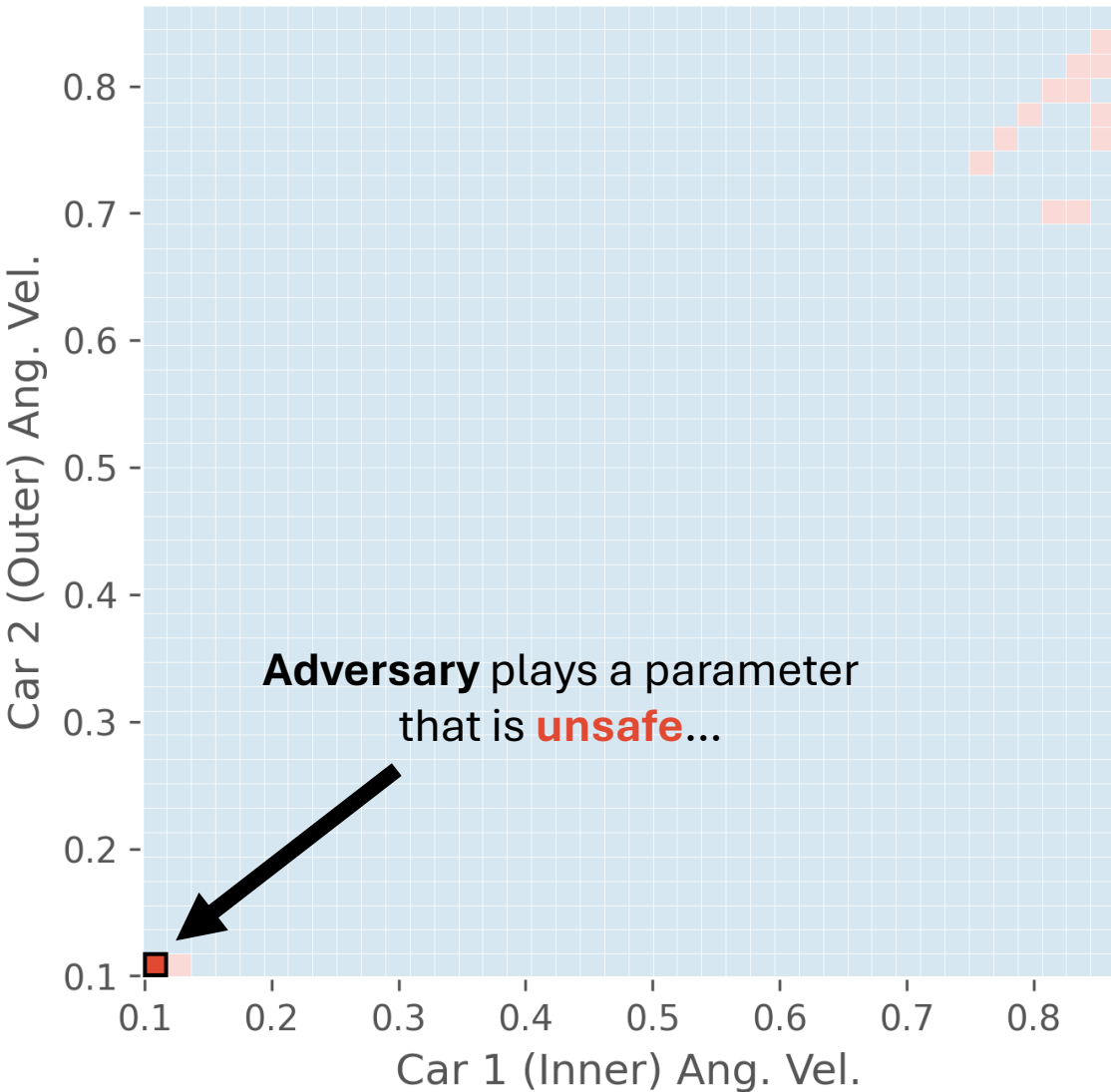
Car 2 (Outer) Ang. Vel.



How do we train a policy that is **more robust** against environment parameters?

Adversarial RL helps but has flaws

Unsafe Safe

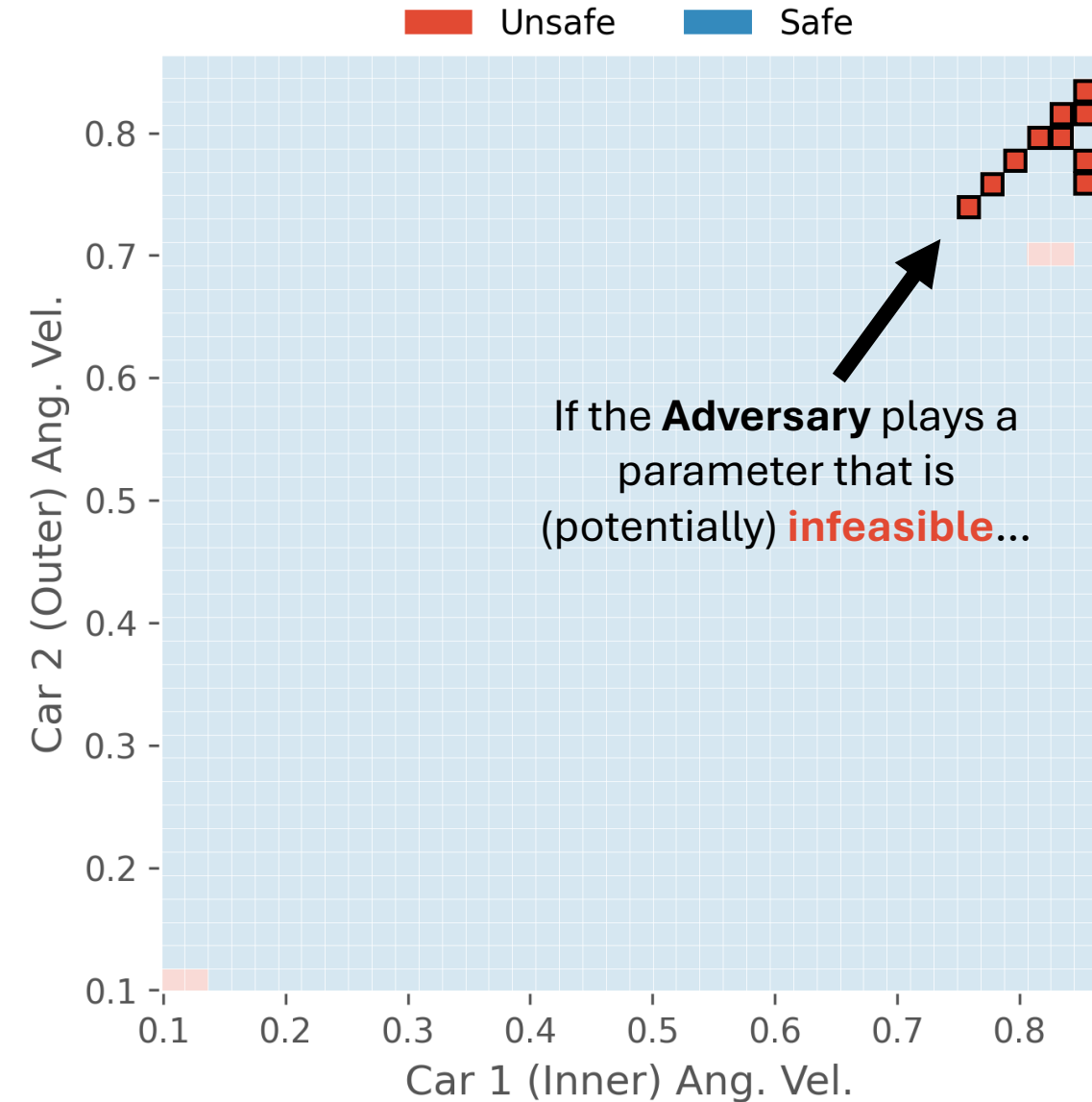


Finetune

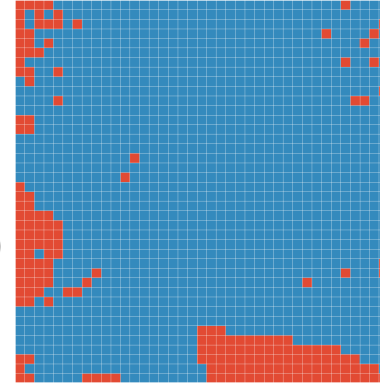
Unsafe Safe



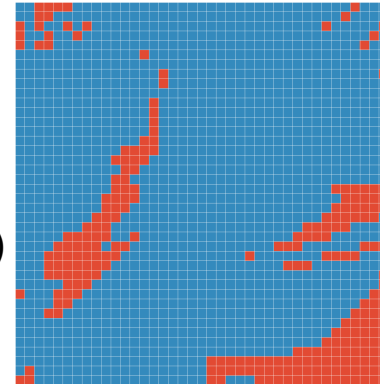
Adversarial RL helps but has flaws



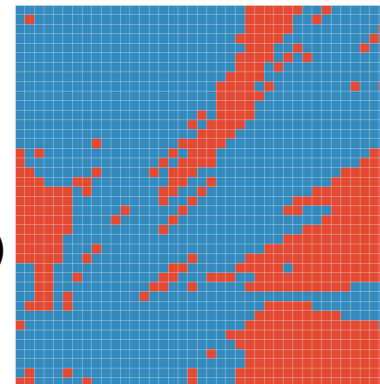
Step 1000
(88% safe)



Step 2000
(86% safe)



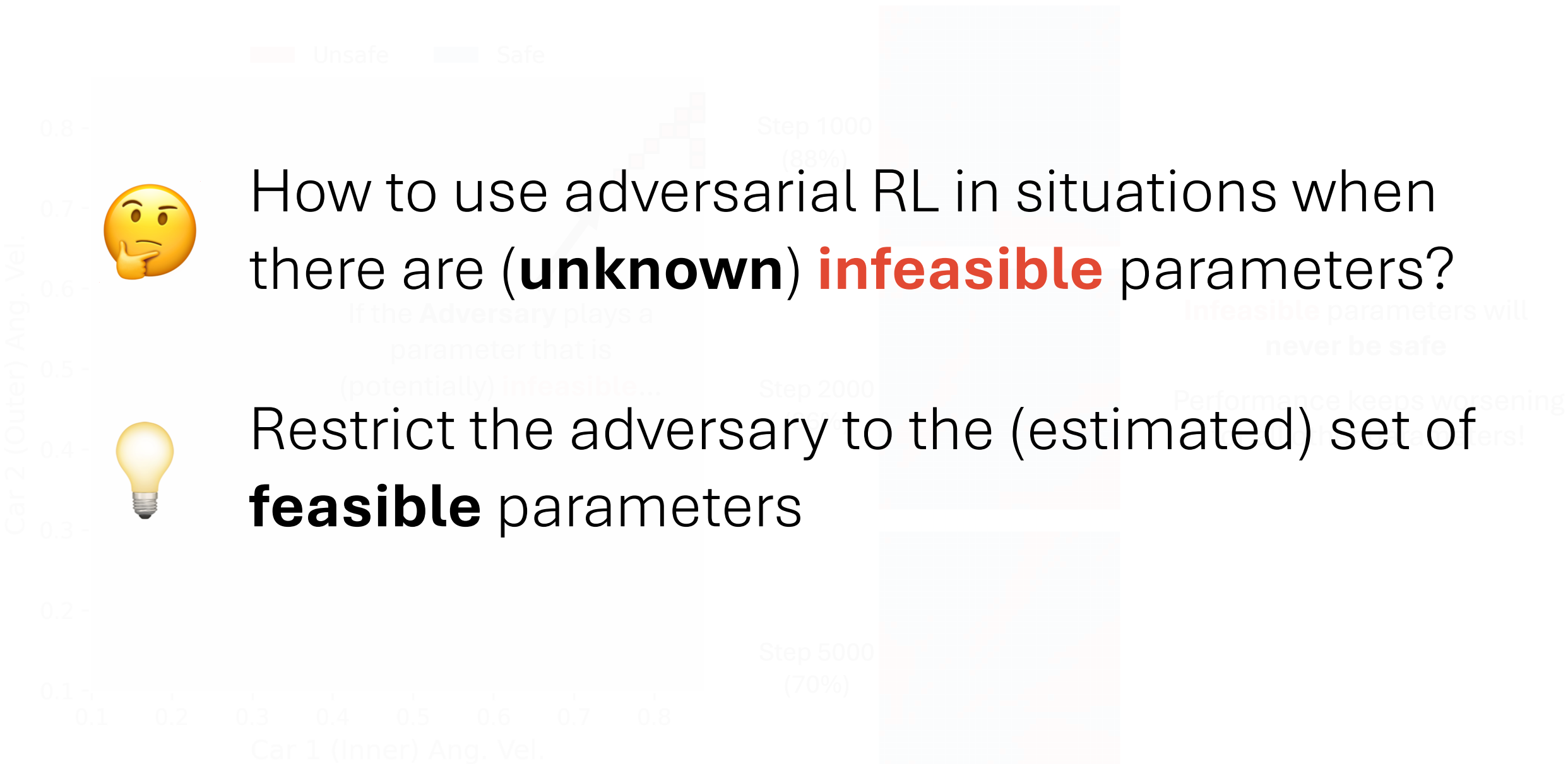
Step 5000
(70% safe)



Infeasible parameters will **never be safe**

Performance keeps worsening for all other parameters!

Adversarial RL helps but has flaws



How to use adversarial RL in situations when there are (**unknown**) **infeasible** parameters?



Restrict the adversary to the (estimated) set of **feasible** parameters

Formal Problem Definition

$$\max_{\Theta' \subseteq \Theta} |\Theta'| \quad \leftarrow \text{Find the largest set of parameters...}$$

$$\text{s.t. } \exists \pi_{(\cdot)} \in \mathcal{M}, \forall \theta \in \Theta',$$

$$\mathbb{1}\{\text{initial state } s_0(\theta) \text{ is safe under } \pi_\theta\} > 0$$

... such that there exists a policy that makes every parameter safe \odot

$$\odot = \max_{\pi_{(\cdot)} \in \mathcal{M}} \min_{\theta \in \Theta'} \sum_{k=0}^{T_{\theta, s_0}} \mathbb{1}\{\mathbf{s}_k \text{ is safe}\} > 0,$$

T_{θ, s_0} \leftarrow Terminate on unsafe state

$$\text{s.t. } \mathbf{s}_0 = s_0(\theta),$$

$$\mathbf{s}_{k+1} = f_\theta(\mathbf{s}_k, \pi_\theta(\mathbf{s}_k))$$

This is a **minimax** RL problem

Building up our method

1. How do we solve for a robust policy, assuming we know the feasible set (of parameters)?
2. How do we (conservatively) estimate the feasible set?
3. How do we improve our feasible set estimate?

Building up our method

1. How do we solve for a robust policy, assuming we know the feasible set (of parameters)?

 Saddle-point finding using techniques from online learning

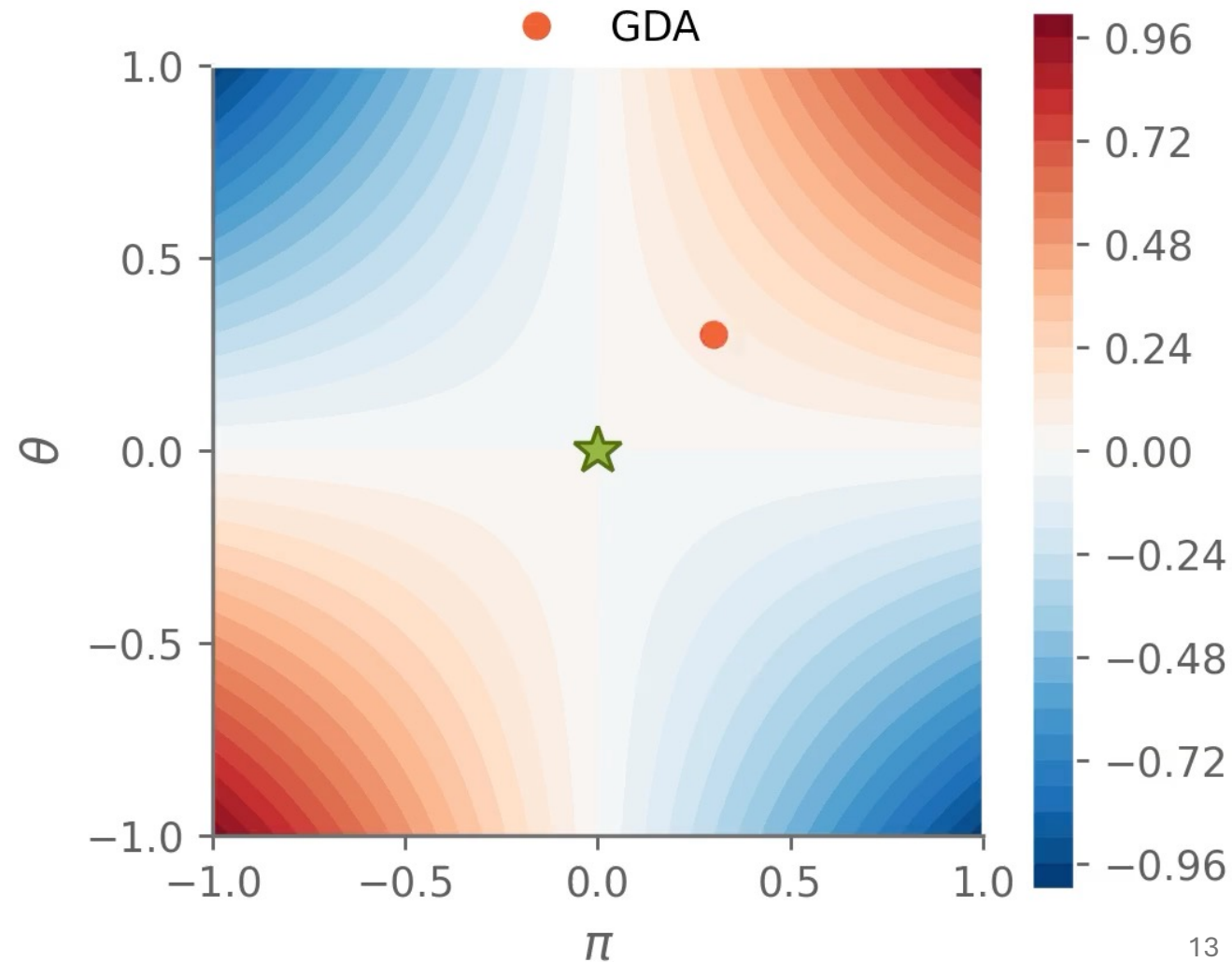
2. How do we (conservatively) estimate the feasible set?

3. How do we improve our feasible set estimate?

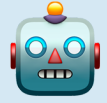
Saddle-point finding

$$\max_{\pi \in \mathcal{M}} \min_{\theta \in \Theta^*} J(\pi, \theta)$$

$$\max_{\pi \in [-1, 1]} \min_{\theta \in [-1, 1]} \pi \cdot \theta$$



Saddle-point finding



Protagonist (FTRL)

$$\pi_{t+1} := \arg \max_{\pi} \underbrace{-\psi_t(\pi) + \sum_{i=1}^t J(\pi, \theta_i)}_{:= \tilde{J}_t(\pi)}$$

“Trust-region” regularization Consider all previous θ_i

θ_{t+1} ↗

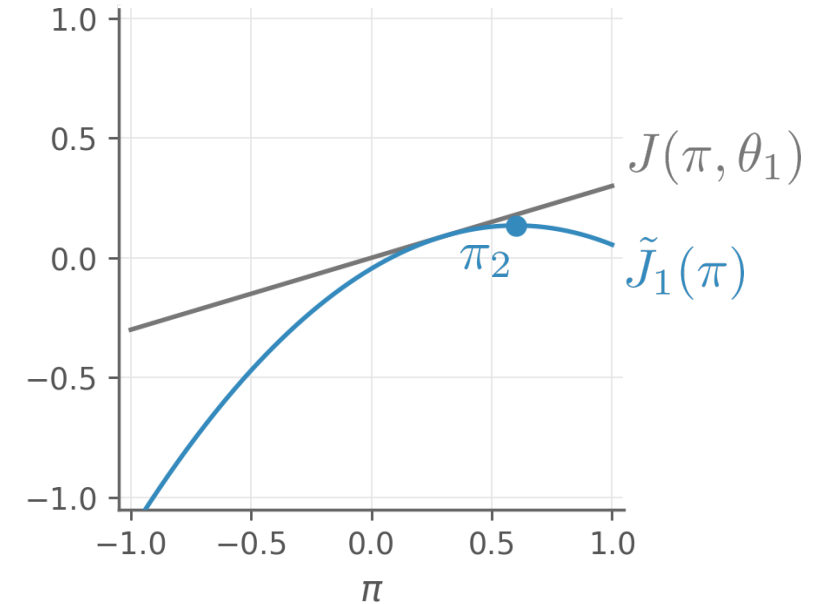
↘ π_{t+1}



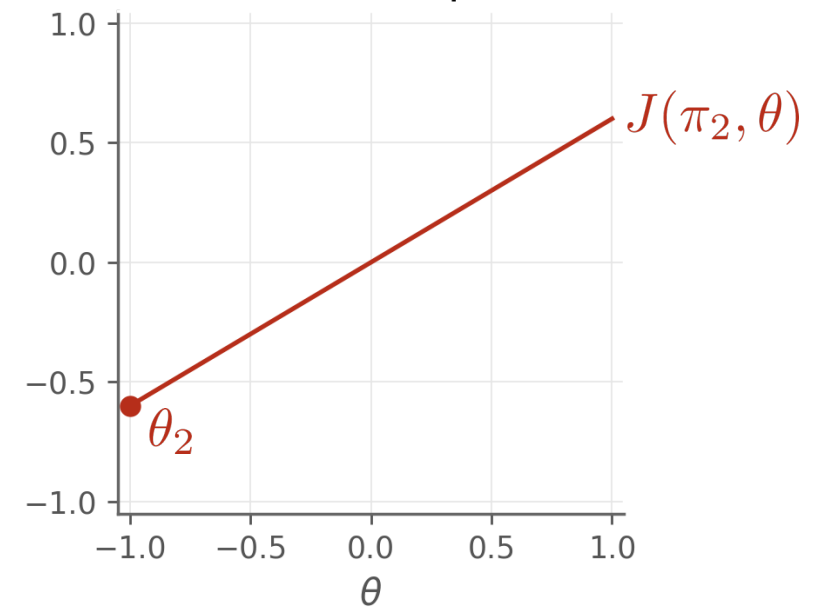
Adversary (Best-Response)

$$\theta_{t+1} := \arg \min_{\theta \in \Theta^*} J(\pi_{t+1}, \theta)$$

FGE π , step 1



FGE θ , step 1



Saddle-point finding

Protagonist (FTRL)

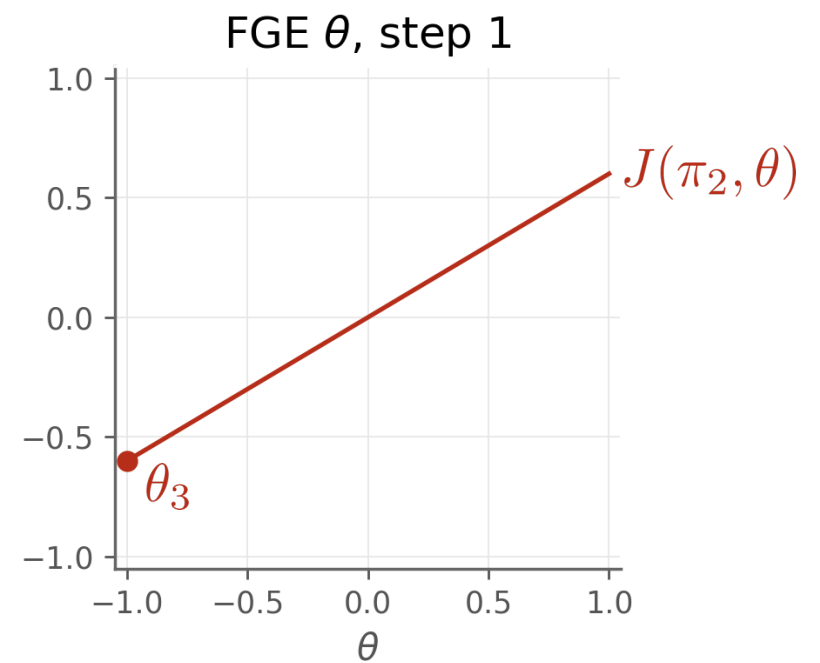
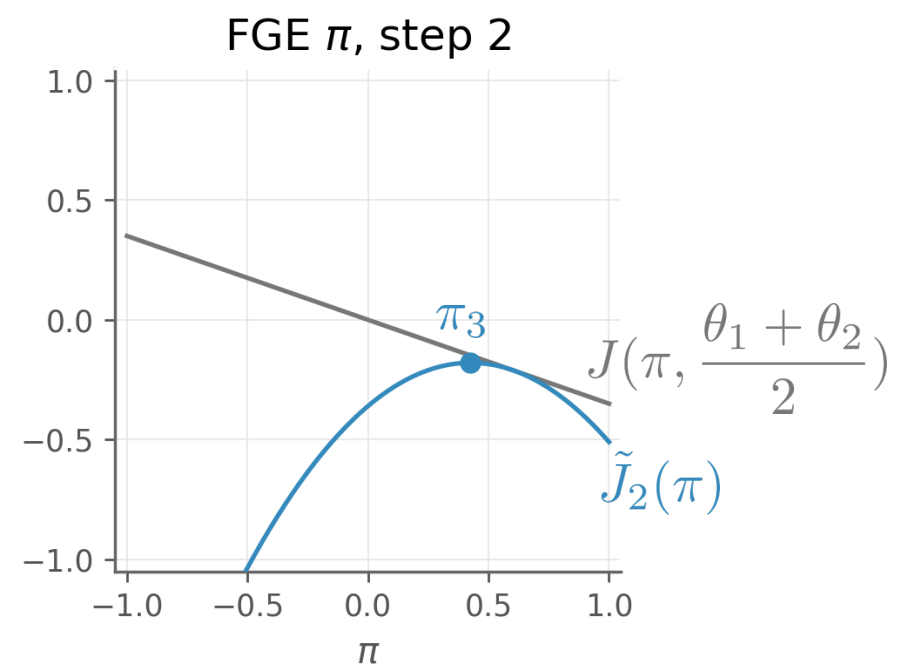
$$\pi_{t+1} := \arg \max_{\pi} \underbrace{-\psi_t(\pi) + \sum_{i=1}^t J(\pi, \theta_i)}_{:= \tilde{J}_t(\pi)}$$

“Trust-region” regularization Consider all previous θ_i



Adversary (Best-Response)

$$\theta_{t+1} := \arg \min_{\theta \in \Theta^*} J(\pi_{t+1}, \theta)$$



Saddle-point finding



Protagonist (FTRL)

$$\pi_{t+1} := \arg \max_{\pi} \underbrace{-\psi_t(\pi) + \sum_{i=1}^t J(\pi, \theta_i)}_{:= \tilde{J}_t(\pi)}$$

“Trust-region” regularization Consider **all** previous θ_i

θ_{t+1} ↗

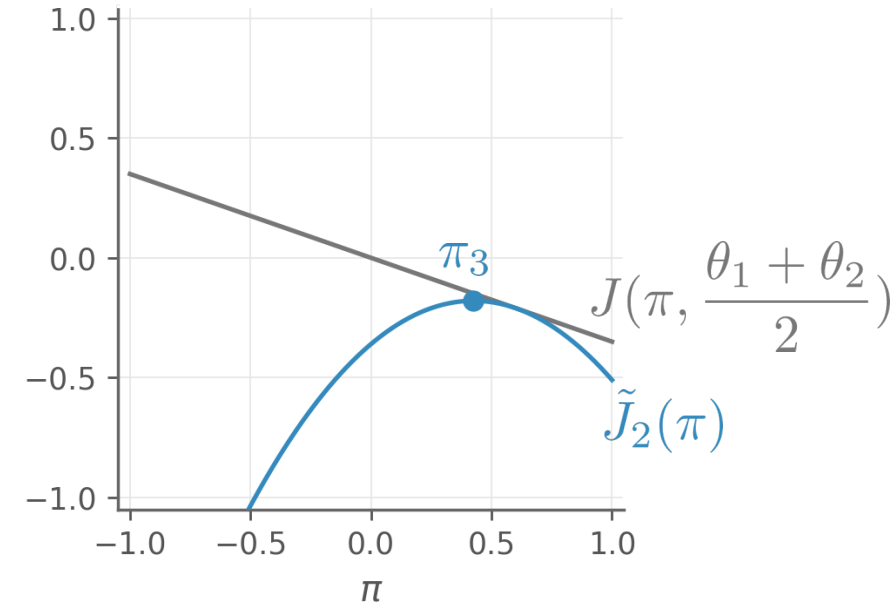
↘ π_{t+1}



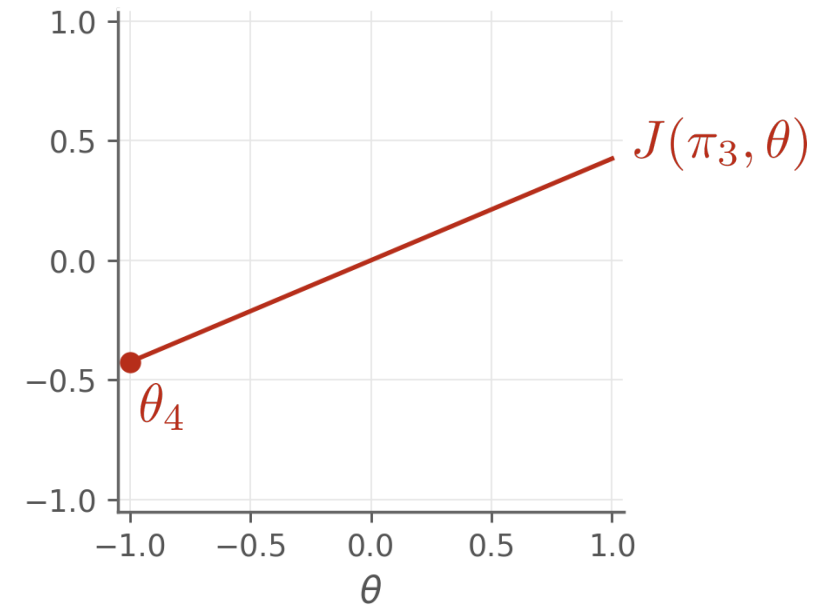
Adversary (Best-Response)

$$\theta_{t+1} := \arg \min_{\theta \in \Theta^*} J(\pi_{t+1}, \theta)$$

FGE π , step 2



FGE θ , step 2



Saddle-point finding



Protagonist (FTRL)

$$\pi_{t+1} := \arg \max_{\pi} \underbrace{-\psi_t(\pi) + \sum_{i=1}^t J(\pi, \theta_i)}_{:= \tilde{J}_t(\pi)}$$

“Trust-region” regularization Consider **all** previous θ_i

θ_{t+1} ↗

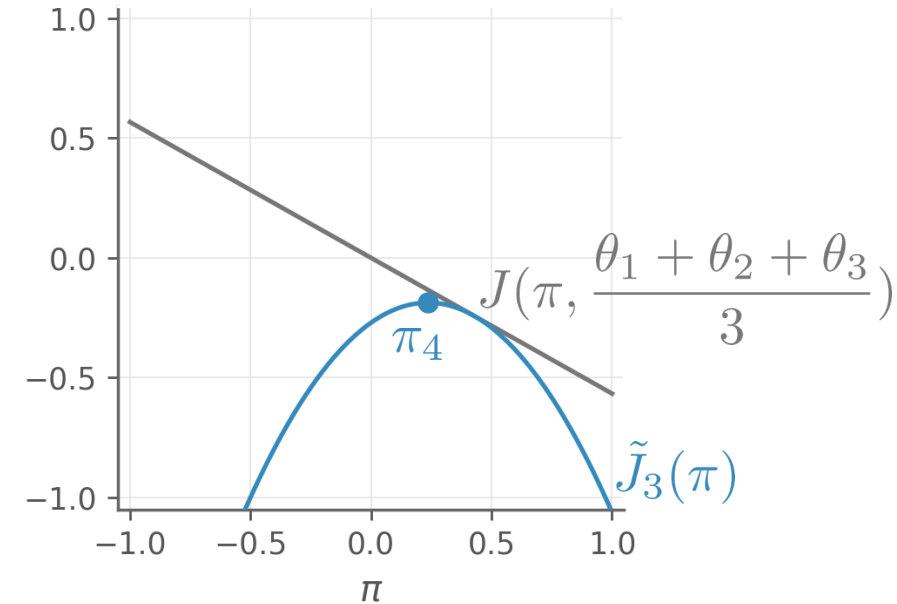
↘ π_{t+1}



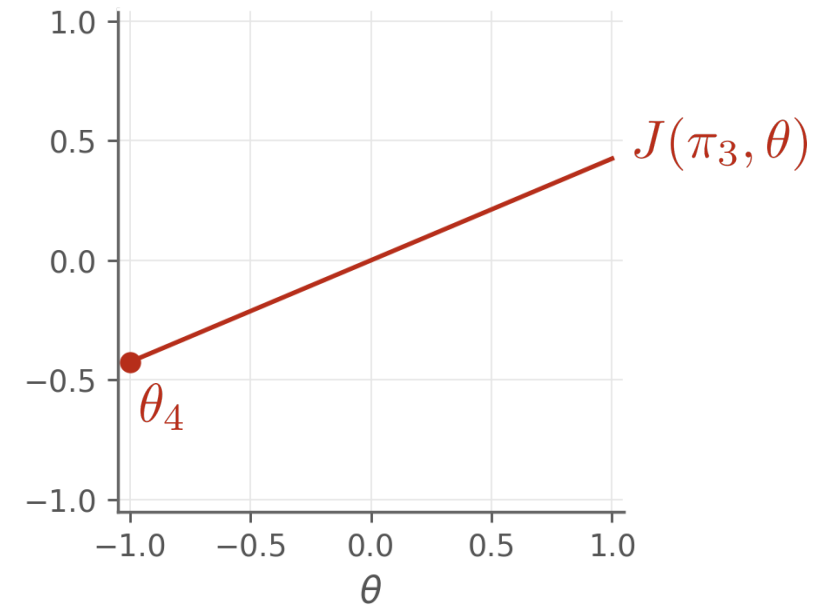
Adversary (Best-Response)

$$\theta_{t+1} := \arg \min_{\theta \in \Theta^*} J(\pi_{t+1}, \theta)$$

FGE π , step 3



FGE θ , step 2



Saddle-point finding

Protagonist (FTRL)

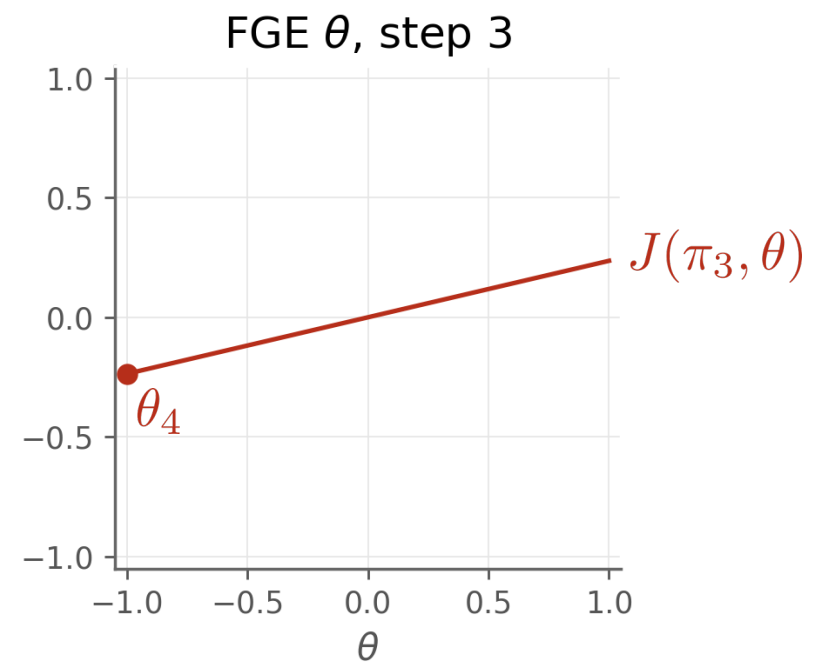
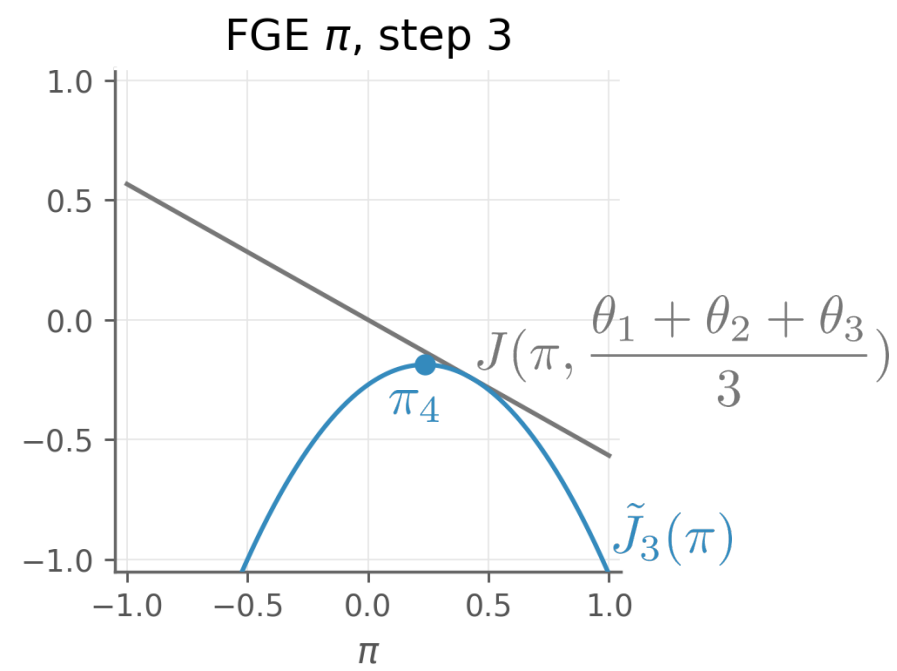
$$\pi_{t+1} := \arg \max_{\pi} \underbrace{-\psi_t(\pi) + \sum_{i=1}^t J(\pi, \theta_i)}_{:= \tilde{J}_t(\pi)}$$

“Trust-region” regularization Consider **all** previous θ_i



Adversary (Best-Response)

$$\theta_{t+1} := \arg \min_{\theta \in \Theta^*} J(\pi_{t+1}, \theta)$$



Saddle-point finding



Protagonist (FTRL)

$$\pi_{t+1} := \arg \max_{\pi} \underbrace{-\psi_t(\pi) + \sum_{i=1}^t J(\pi, \theta_i)}_{:= \tilde{J}_t(\pi)}$$

“Trust-region” regularization Consider **all** previous θ_i

θ_{t+1} ↗

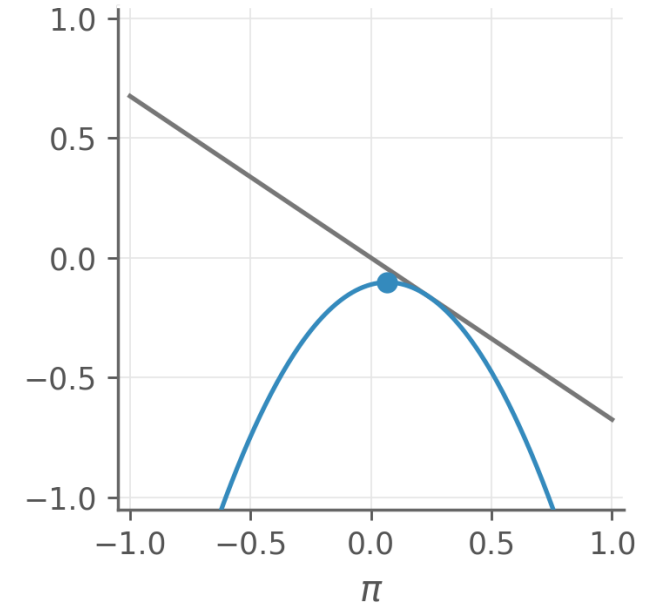
↘ π_{t+1}



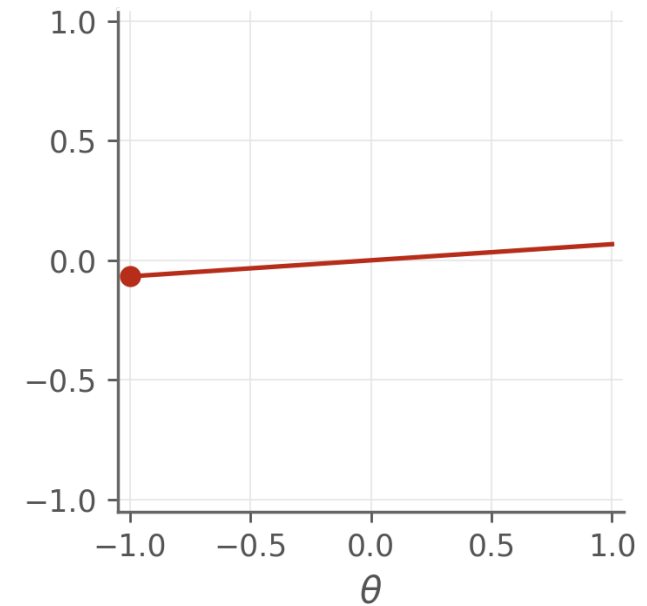
Adversary (Best-Response)

$$\theta_{t+1} := \arg \min_{\theta \in \Theta^*} J(\pi_{t+1}, \theta)$$

FGE π , step 4



FGE θ , step 4



Saddle-point finding



Protagonist (FTRL)

$$\pi_{t+1} := \arg \max_{\pi} \underbrace{-\psi_t(\pi) + \sum_{i=1}^t J(\pi, \theta_i)}_{:= \tilde{J}_t(\pi)}$$

“Trust-region” regularization Consider **all** previous θ_i

θ_{t+1} ↗

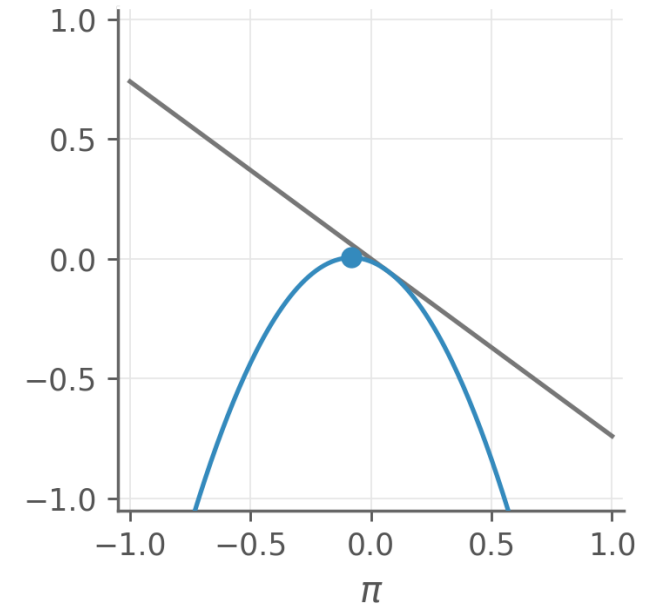
↘ π_{t+1}



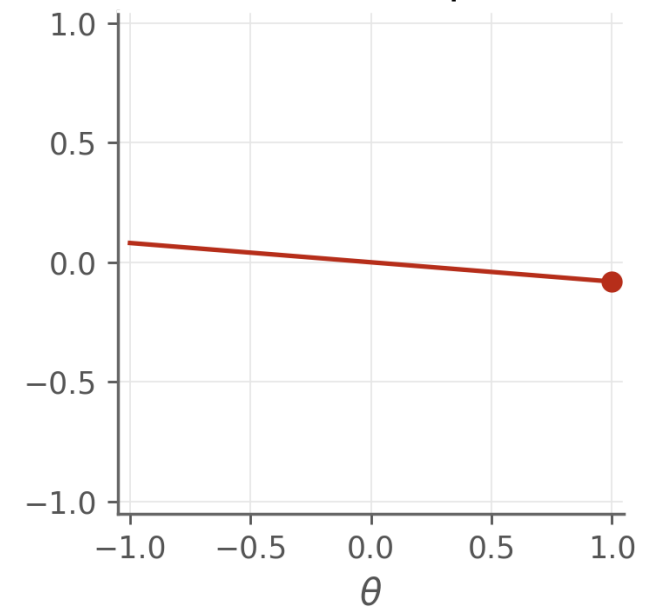
Adversary (Best-Response)

$$\theta_{t+1} := \arg \min_{\theta \in \Theta^*} J(\pi_{t+1}, \theta)$$

FGE π , step 5



FGE θ , step 5



Saddle-point finding



Protagonist (FTRL)

$$\pi_{t+1} := \arg \max_{\pi} \underbrace{-\psi_t(\pi) + \sum_{i=1}^t J(\pi, \theta_i)}_{:= \tilde{J}_t(\pi)}$$

“Trust-region” regularization Consider **all** previous θ_i

θ_{t+1} ↗

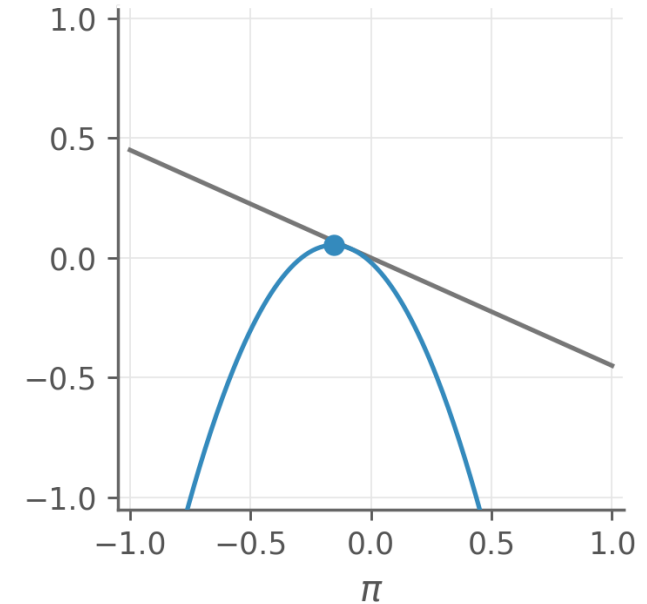
↘ π_{t+1}



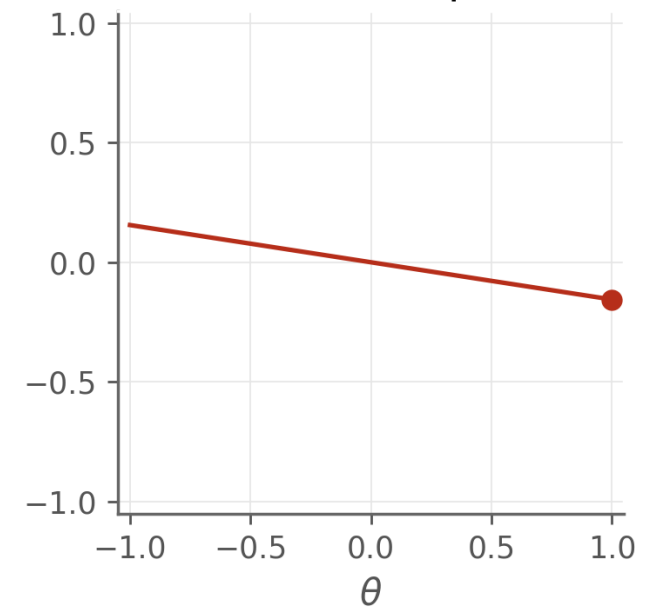
Adversary (Best-Response)

$$\theta_{t+1} := \arg \min_{\theta \in \Theta^*} J(\pi_{t+1}, \theta)$$

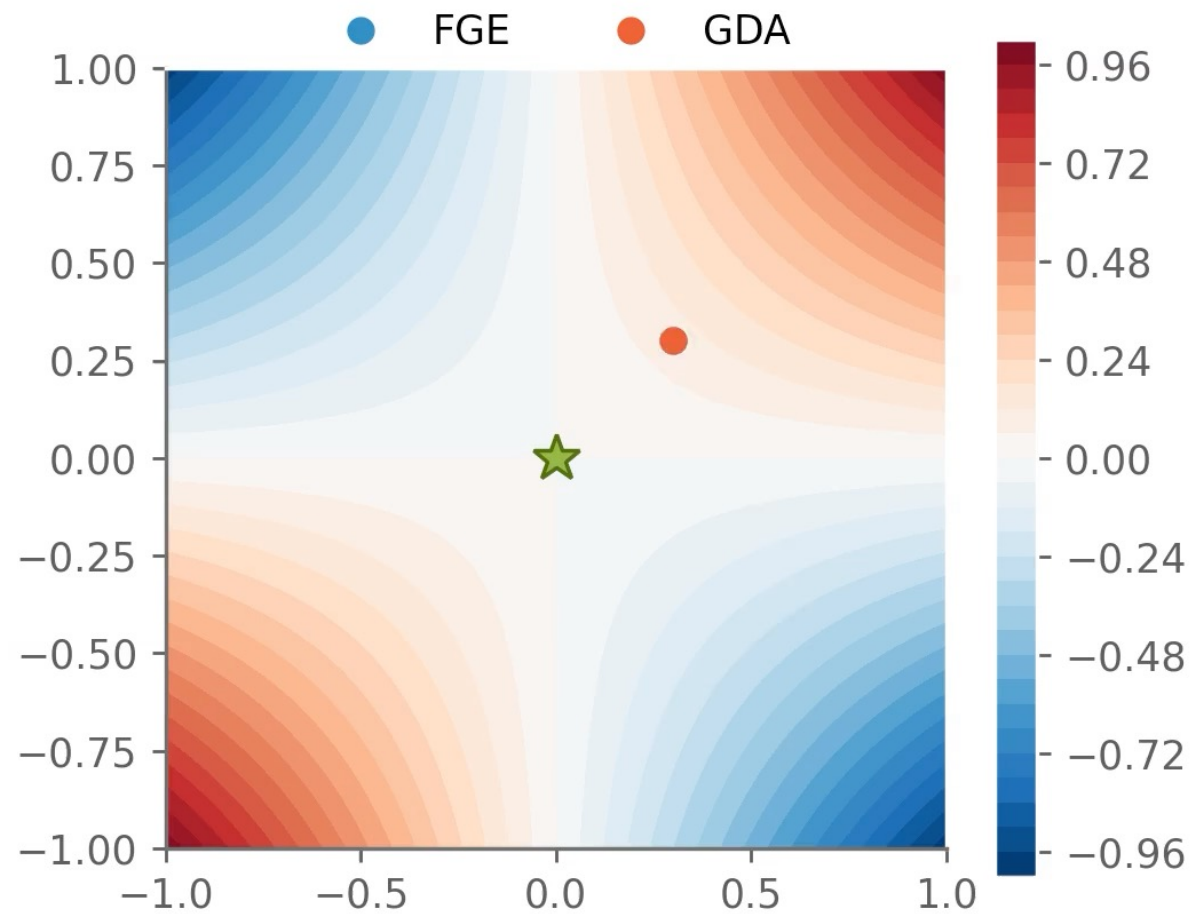
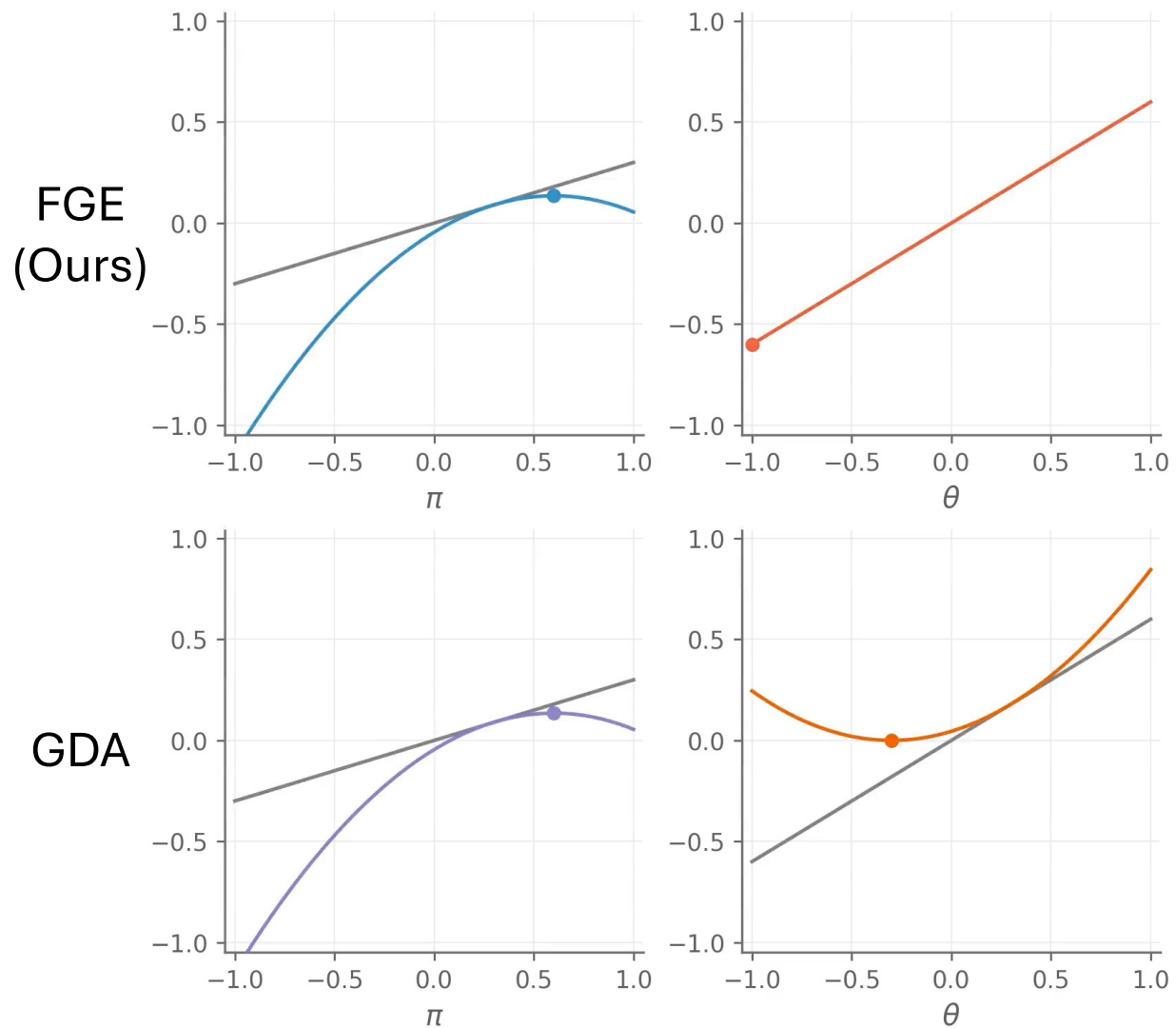
FGE π , step 6



FGE θ , step 6



Saddle-point finding



Saddle-point finding | Implementation



Protagonist (FTRL)

$$\pi_{t+1} := \arg \max_{\pi} \underbrace{-\psi_t(\pi) + \sum_{i=1}^t J(\pi, \theta_i)}_{:= \tilde{J}_t(\pi)}$$

$$\pi_{t+1} = \pi_t + \eta_t \nabla_{\pi} \frac{1}{t} \sum_{i=1}^t J(\pi, \theta_i)$$

This is just **gradient ascent!**
→ PPO update with **domain randomization** over θ

$$= \pi_t + \eta_t \nabla_{\pi} \mathbb{E}_{\theta \sim \mathcal{D}_{\theta,t}} [J(\pi, \theta)], \quad \mathcal{D}_{\theta,t} = \text{Uniform}(\{\theta_i\}_{i=1}^t)$$

Saddle-point finding | Implementation

Robust RL

On-Policy RL

PPO

Building up our method

1. How do we solve for a robust policy, assuming we know the feasible set (of parameters)?

💡 Saddle-point finding using techniques from online learning

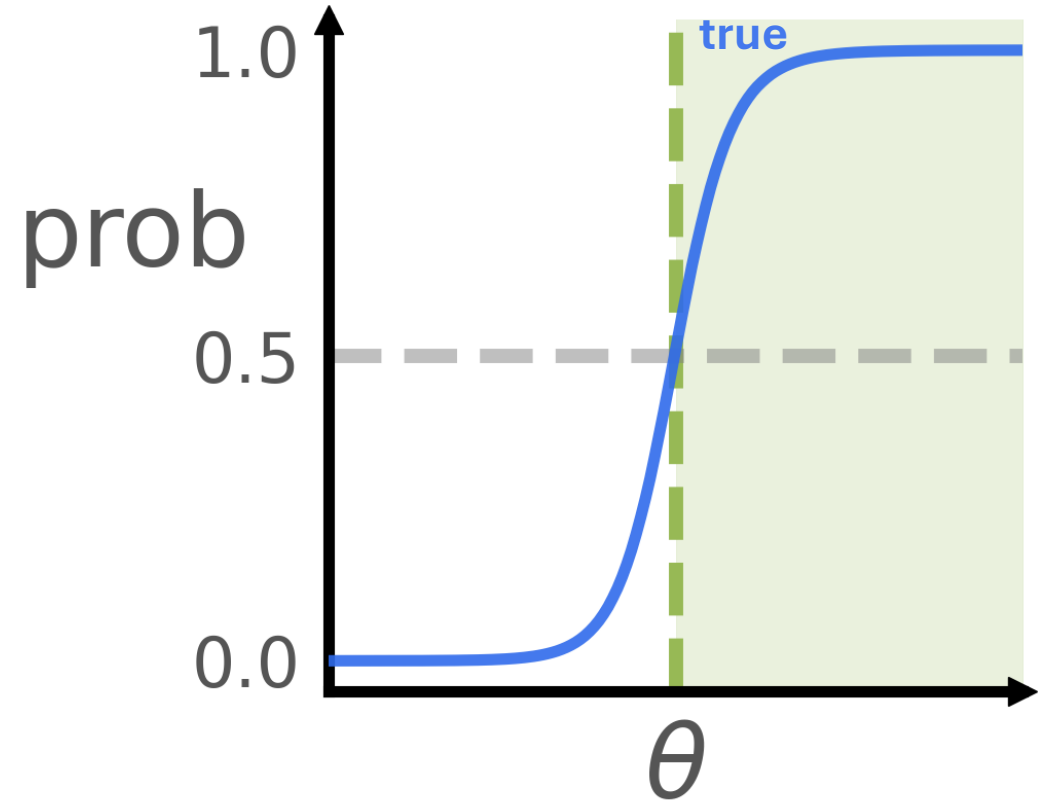
2. How do we (conservatively) estimate the feasible set?

💡 Use policy rollouts in a smart way

3. How do we improve our feasible set estimate?

Feasible Set Estimation

Want to learn $p(\text{safe} \mid \theta)$.

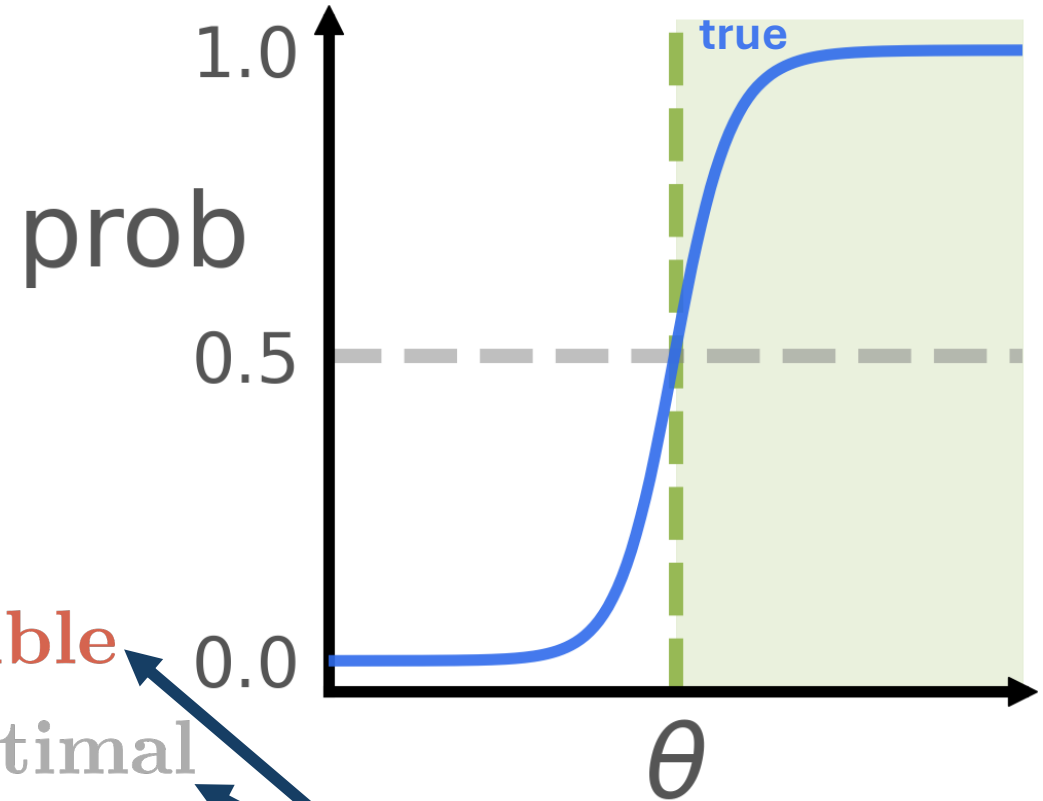


Feasible Set Estimation

Want to learn $p(\text{safe} \mid \theta)$.

π renders θ **safe** \implies θ is **feasible**

unsafe \implies $\left\{ \begin{array}{l} \theta \text{ is } \text{infeasible} \\ \pi \text{ is } \text{suboptimal} \end{array} \right.$



Impossible to distinguish between these two cases!

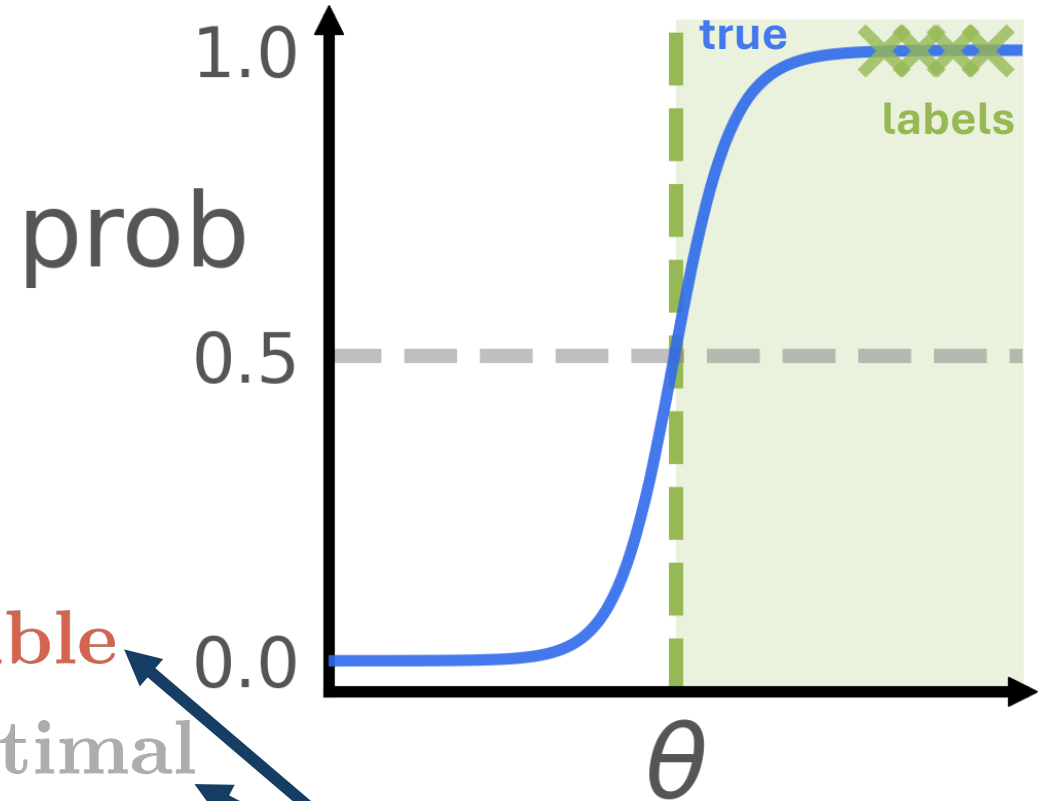
Feasible Set Estimation

Want to learn $p(\text{safe} \mid \theta)$.

π renders θ **safe** \implies θ is **feasible**

unsafe \implies $\left\{ \begin{array}{l} \theta \text{ is } \textbf{infeasible} \\ \pi \text{ is } \textbf{suboptimal} \end{array} \right.$

We have (θ, safe) , but no (θ, unsafe) .



Impossible to distinguish between these two cases!

Feasible Set Estimation

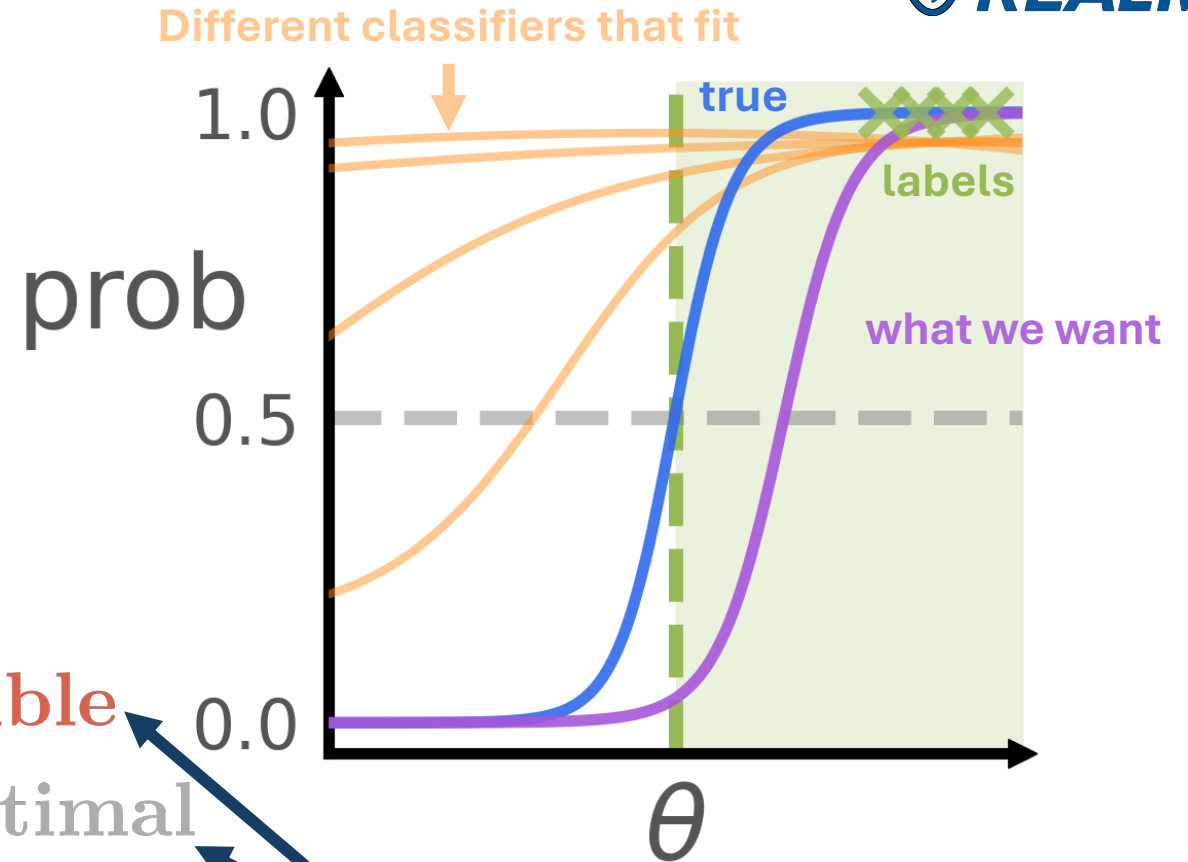
Want to learn $p(\text{safe} \mid \theta)$.

π renders θ **safe** \implies θ is **feasible**

unsafe \implies $\left\{ \begin{array}{l} \theta \text{ is } \textbf{infeasible} \\ \pi \text{ is } \textbf{suboptimal} \end{array} \right.$

We have (θ, safe) , but no (θ, unsafe) .

🌀 Cannot learn a classifier with only positive labels!

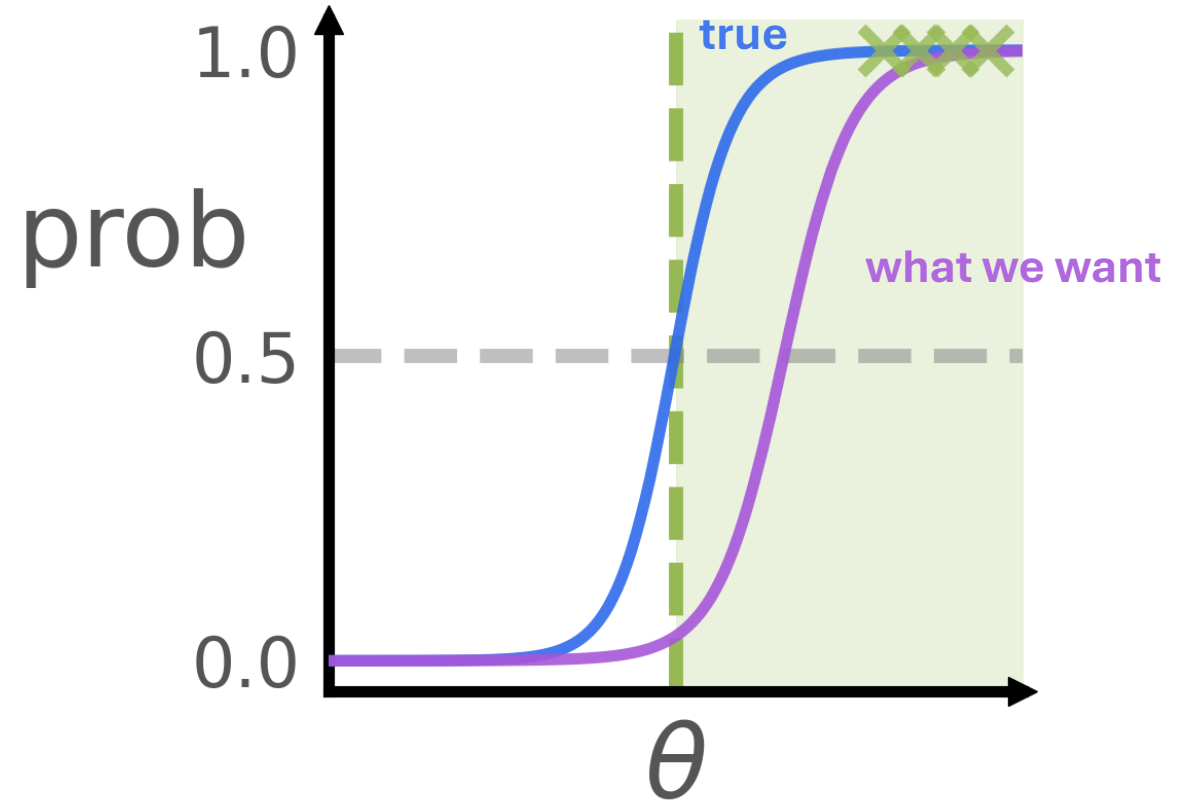


Impossible to distinguish between these two cases!

Feasible Set Estimation

We have (θ, safe) , but no (θ, unsafe) .

🌀 Cannot learn a classifier with only positive labels!



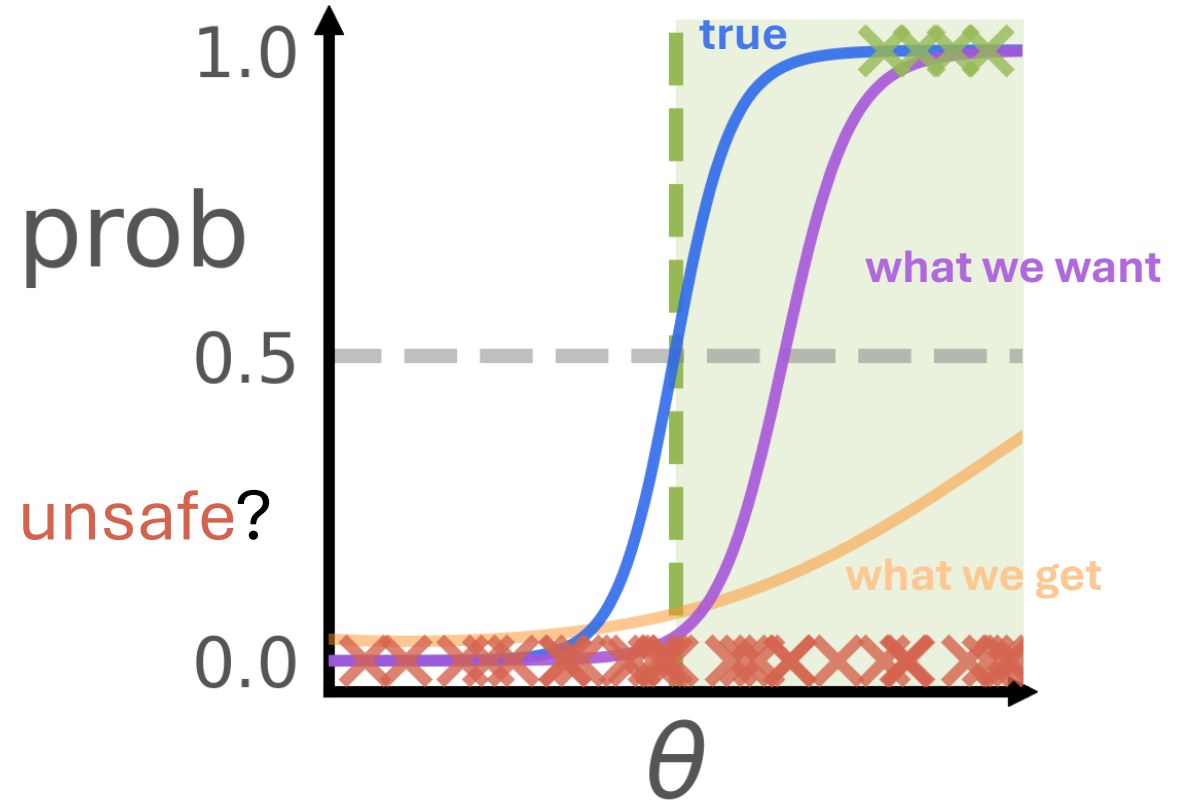
Feasible Set Estimation

We have (θ, safe) , but no (θ, unsafe) .

🌀 Cannot learn a classifier with only positive labels!

🤔 What if we use (θ, unsafe) when π is **unsafe**?

🌀 Poor fit if π degrades!



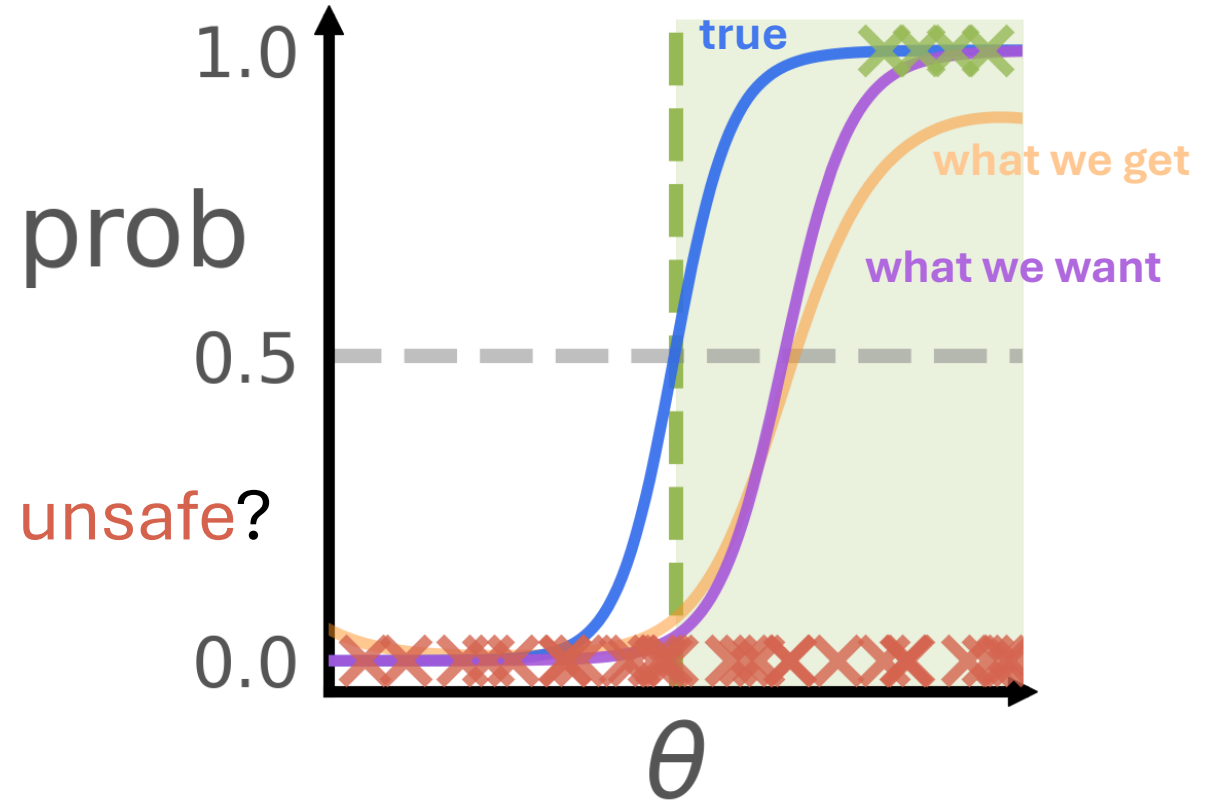
Feasible Set Estimation

We have (θ, safe) , but no (θ, unsafe) .

🌀 Cannot learn a classifier with only positive labels!

🤔 What if we use (θ, unsafe) when π is **unsafe**?

🌀 Poor fit if π degrades!



💡 **Solution:** *Weight* the **safe** samples *more* than the **unsafe** samples!

$$P_{\text{mix}}(\text{safe}, \theta) = \alpha \underbrace{p^*(\text{safe} \mid \theta) p_{\mathcal{D}_{\text{feasible}}}(\theta)}_{\text{observed feasible } \theta} + (1 - \alpha) \underbrace{p^{\pi}(\text{safe} \mid \theta) \rho(\theta)}_{\text{mixed safe/unsafe } \theta}$$

α lets us control the false negative rate!

Feasible Set Estimation | Guarantees

Theorem 1. Let $q_{\psi^*} : \Theta \rightarrow [0, 1]$ be the optimal variational approximation to P_{mix} wrt. variational parameters ψ .

For a **feasible** $\theta \in \mathcal{D}_{\text{feasible}}$, q_{ψ^*} classifies θ as **safe** with probability β where

$$p^\pi(\text{safe} \mid \theta) \geq \beta - (1 - \beta) \frac{\alpha}{1 - \alpha} \frac{p_{\mathcal{D}_{\text{feasible}}}(\theta)}{\rho(\theta)}$$

For an **infeasible** $\theta \in \Theta^{*\mathbb{C}}$, q_{ψ^*} always correctly classifies θ as **unsafe**.

Feasible Set Estimation | Implementation



Feasibility Learning

On-Policy RL

PPO

Building up our method

1. How do we solve for a robust policy, assuming we know the feasible set (of parameters)?

💡 Saddle-point finding using techniques from online learning

2. How do we (conservatively) estimate the feasible set?

💡 Use policy rollouts in a smart way

3. How do we improve our feasible set estimate?

💡 Focus exploration on parameters currently not feasible

Expanding the Feasible Set Estimation

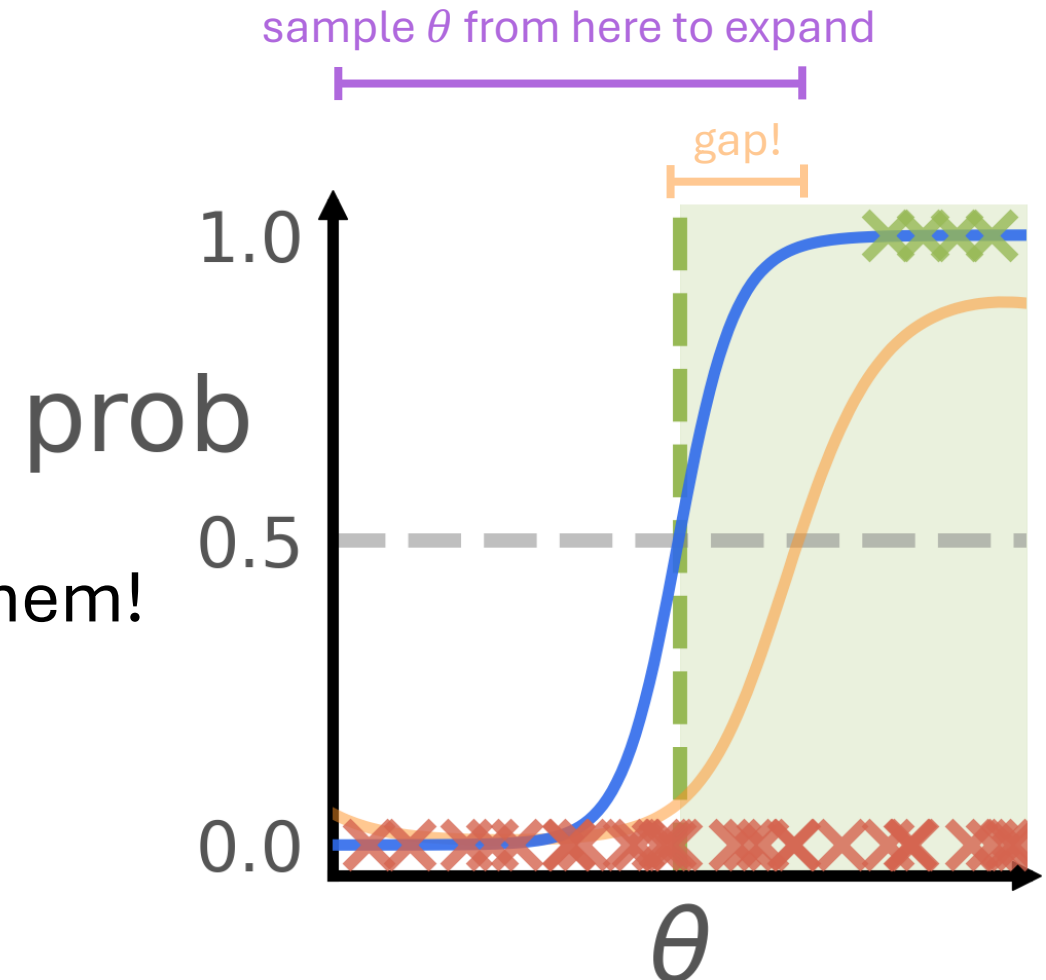
🤔 Estimated feasible set size depends on the quality of the policy.

To expand, need to render **feasible** but currently **unsafe** θ as **safe**.

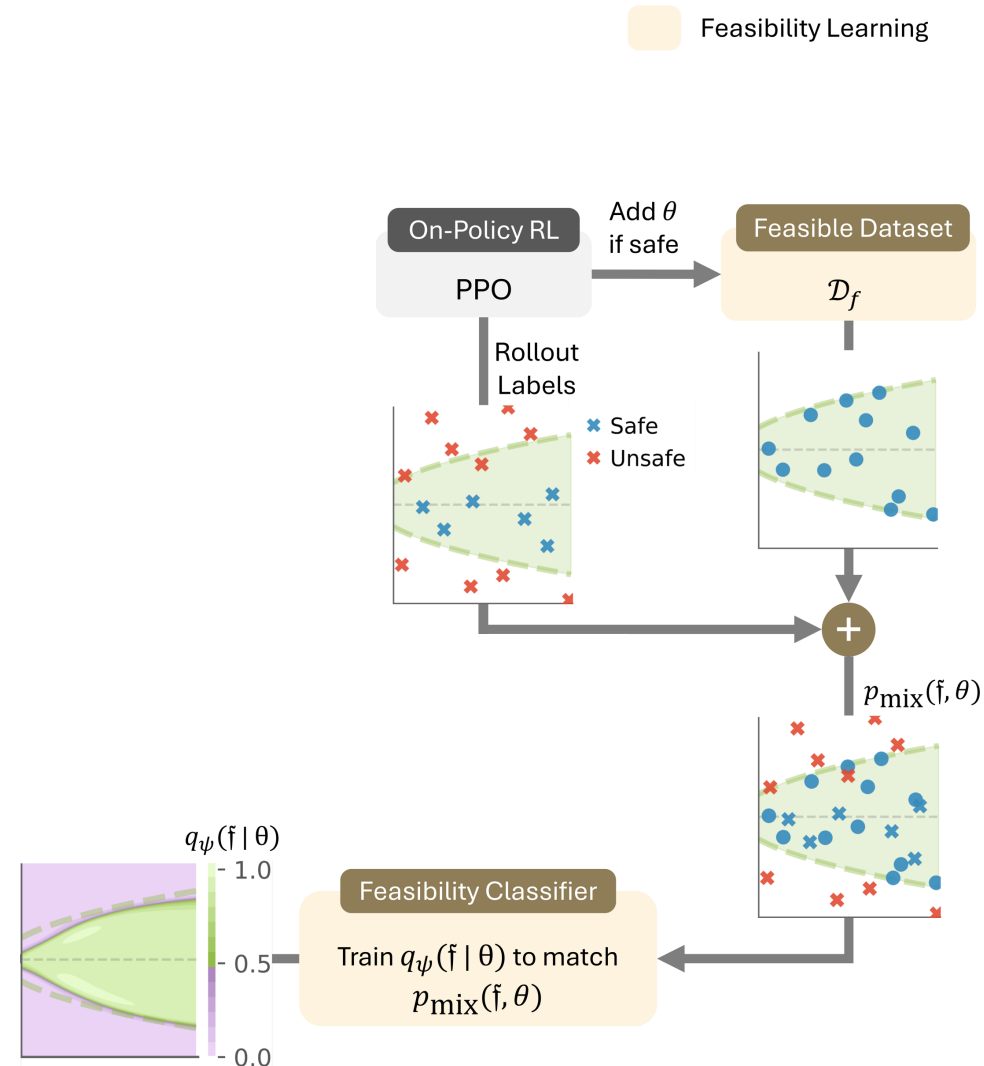
🌀 Which θ are actually feasible?

💡 **Solution**: Rejection sample on all of them!

$$\theta \sim p(\cdot \mid \text{predict } \mathbf{infeasible})$$



Expanding the Feasible Set Estimation Implementation

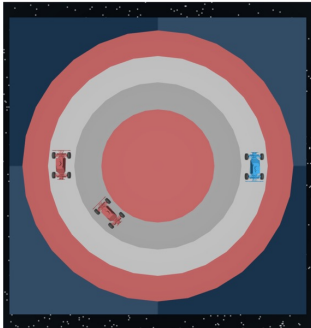


Results | Tasks

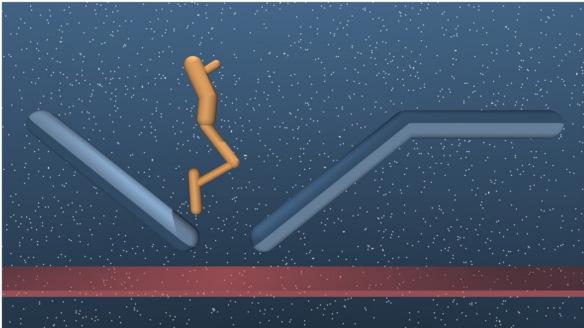
ToyLevels



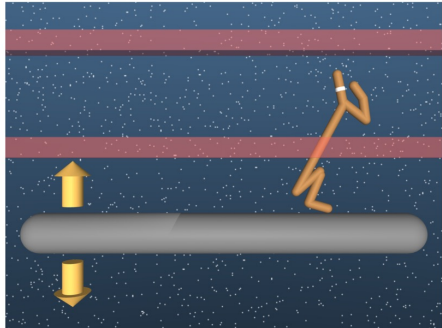
Dubins



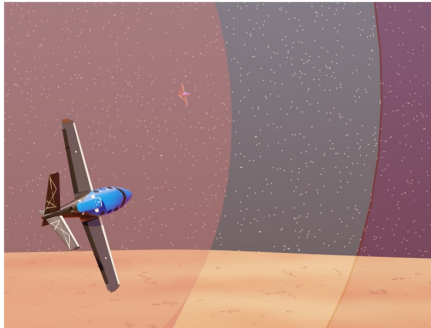
Hopper



HalfCheetah

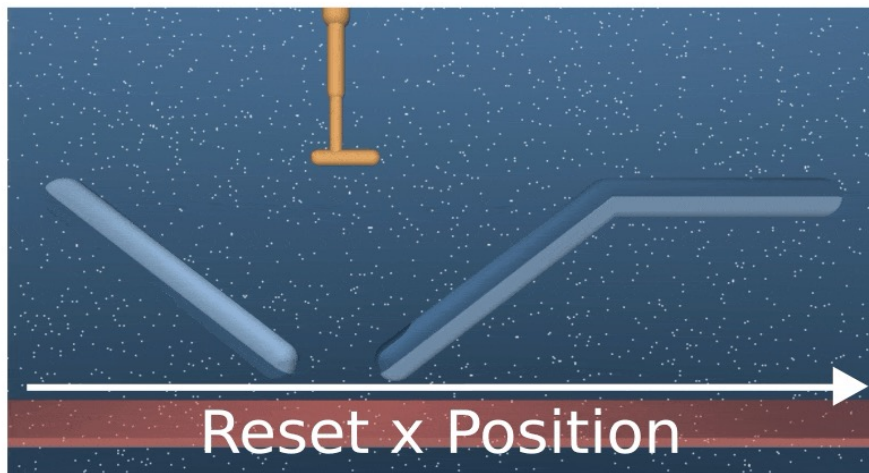
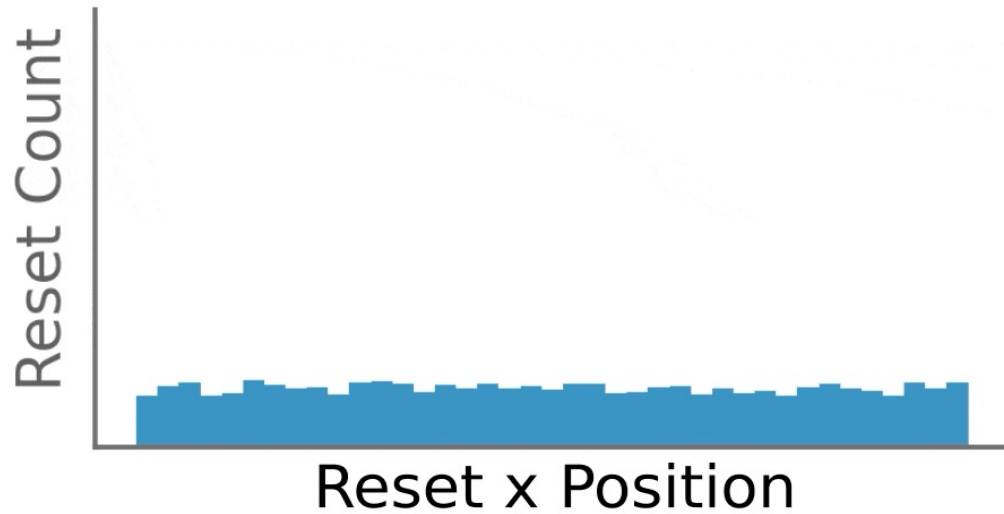


FixedWing

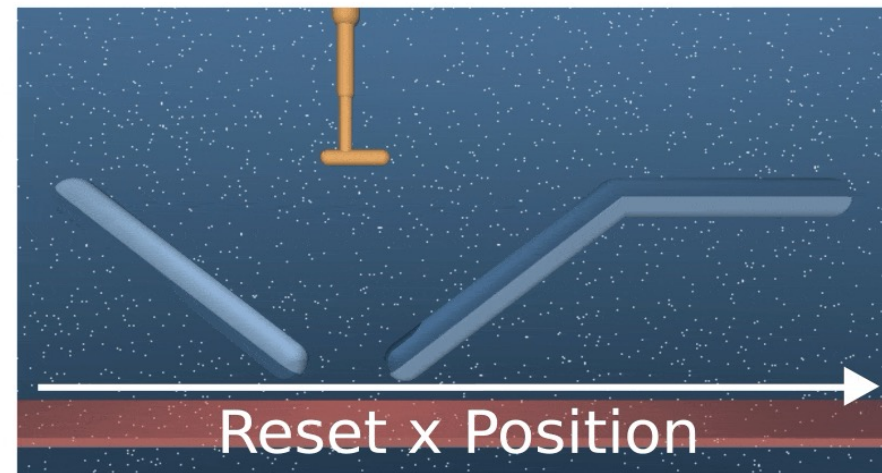
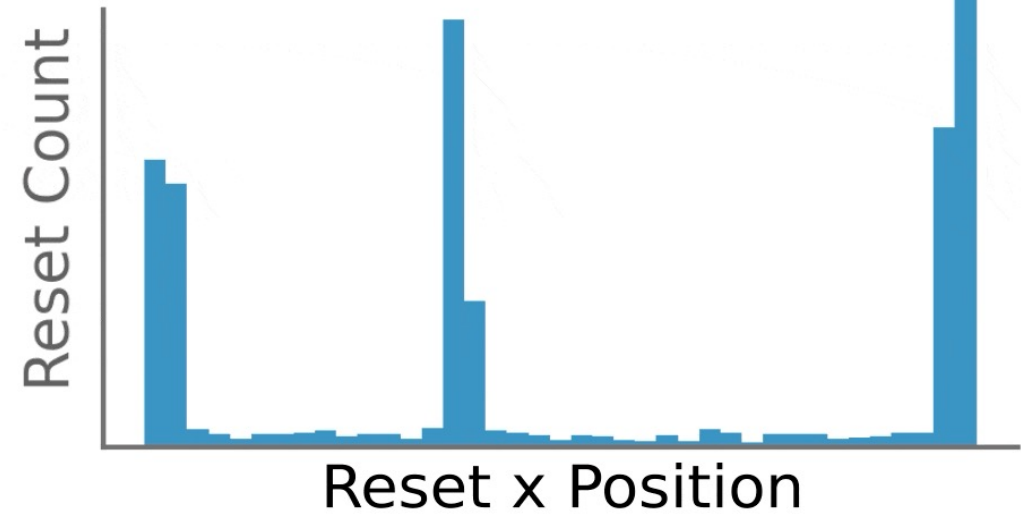


Results | Hopper

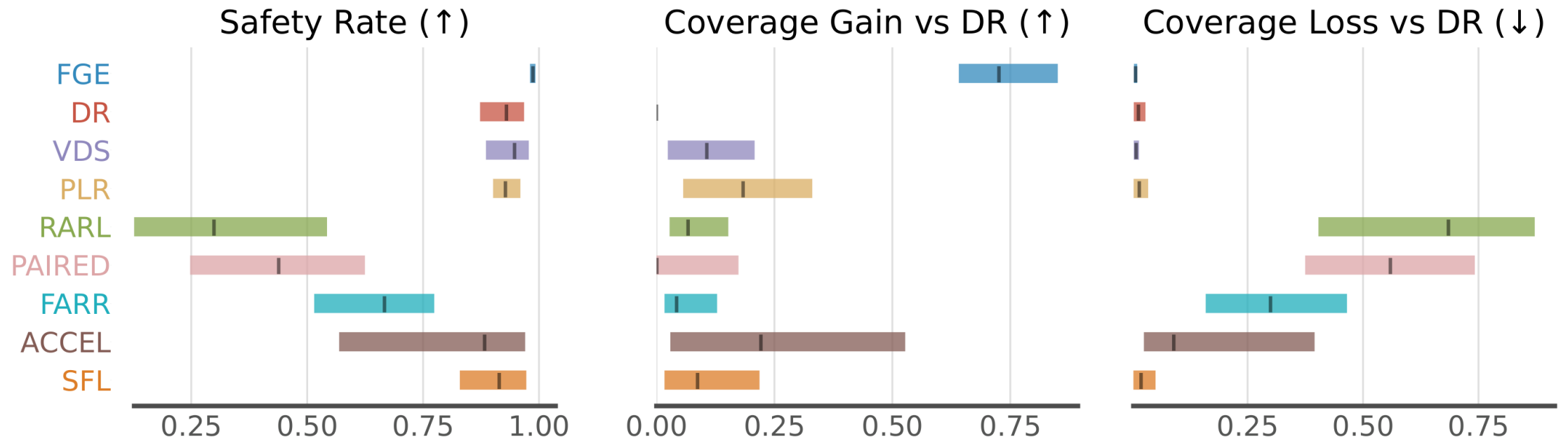
Domain Randomization



FGE (Ours)



Results | Overall Performance



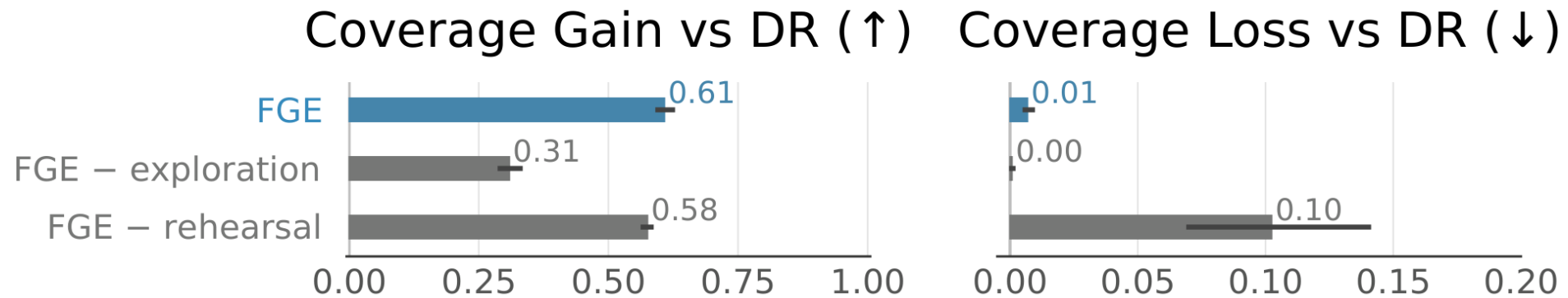
Ablation Studies

(Q1) Is the explore distribution important?

(Q2) Is the rehearsal buffer important?

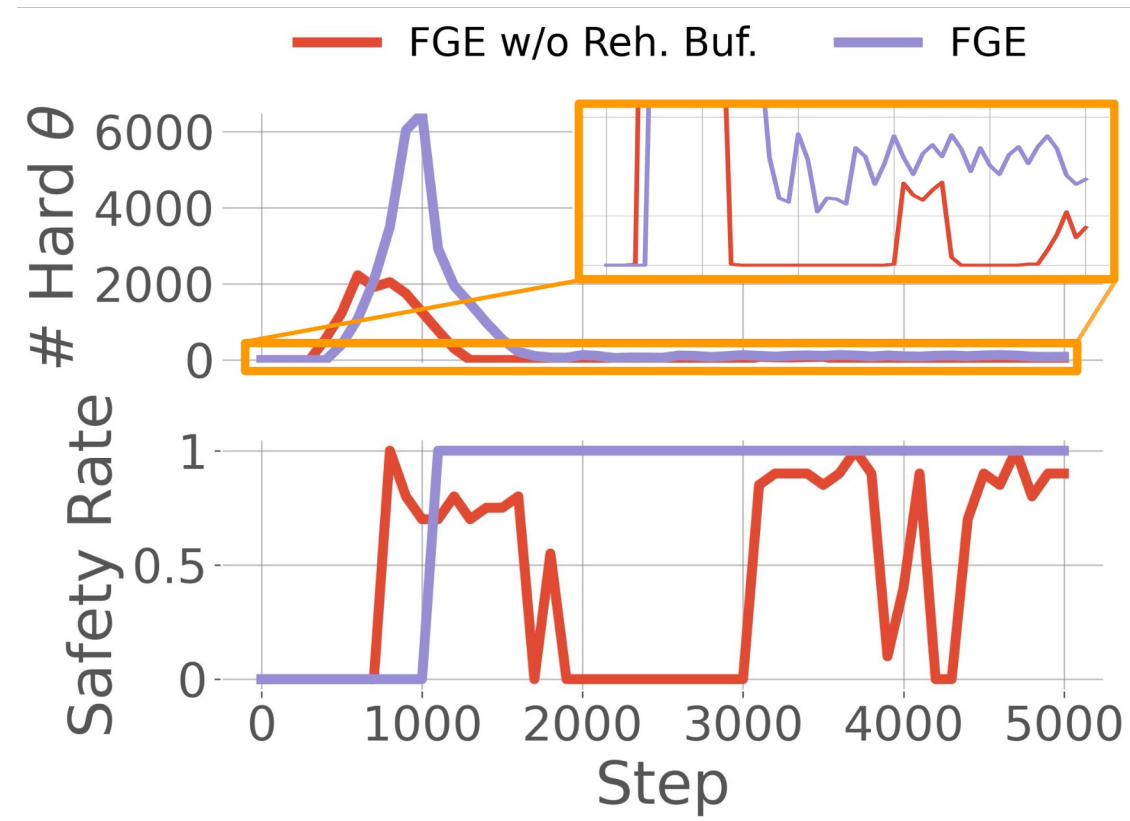
(Q3) Is the feasibility classifier important?

(Q1) Is the explore distribution important?



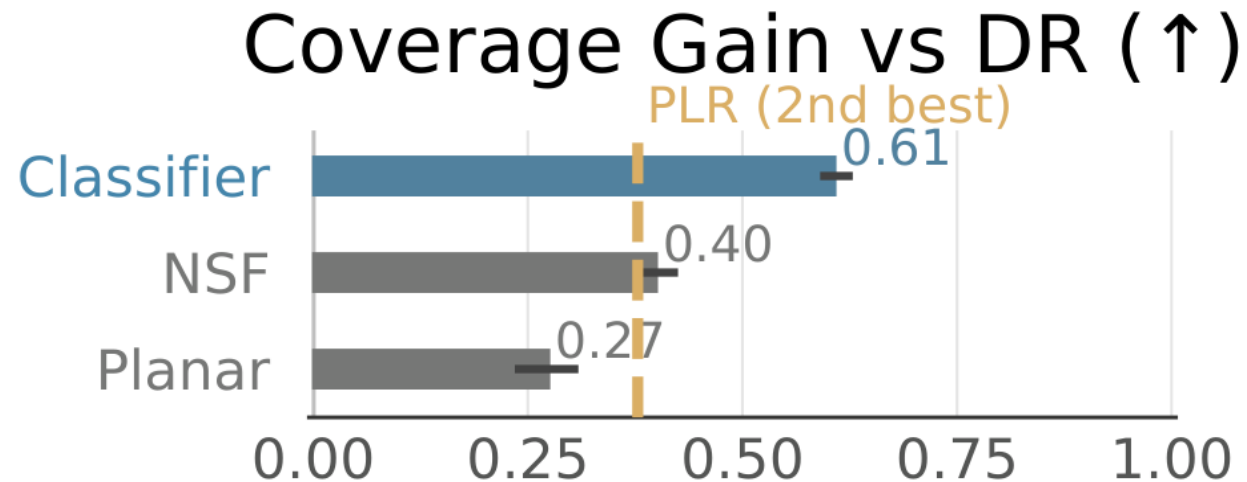
(A) Yes. Removing it leads to smaller feasible sets, and smaller coverage gain vs DR.

(Q2) Is the rehearsal buffer important?



(A) Yes. Removing it leads to catastrophic forgetting on the feasible set.

(Q3) Is the feasibility classifier important?

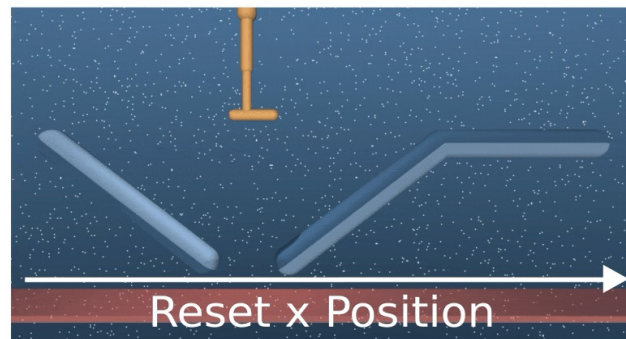
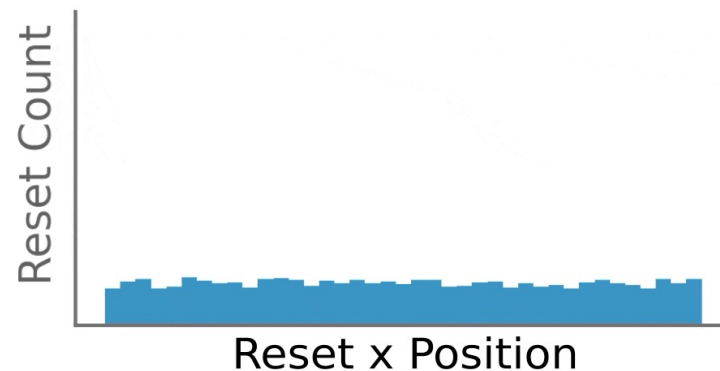


(A) Yes. Using a density model as a classifier instead degrades performance, largely due to thresholding issues.

Summary

- We propose **Feasibility Guided Exploration (FGE)** to improve safety under unknown feasibility using RL.

Domain Randomization



FGE (Ours)

