

Class notes

1. Homework 5 due Tuesday, November 13th 11:59pm

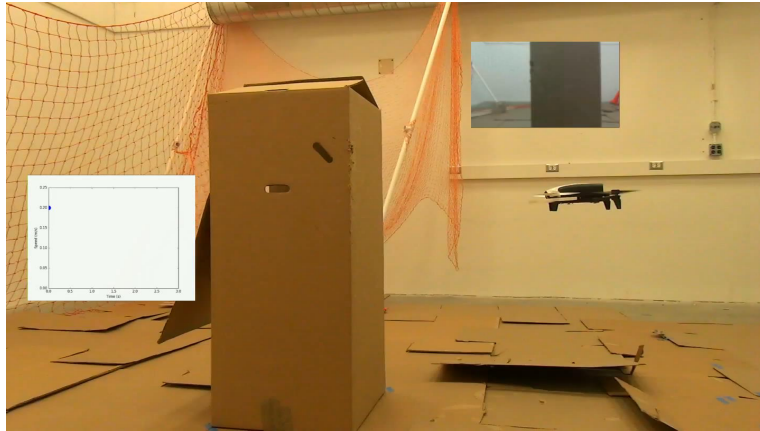
Real-World Robot Learning: Safety and Flexibility

CS294-112: Deep Reinforcement Learning

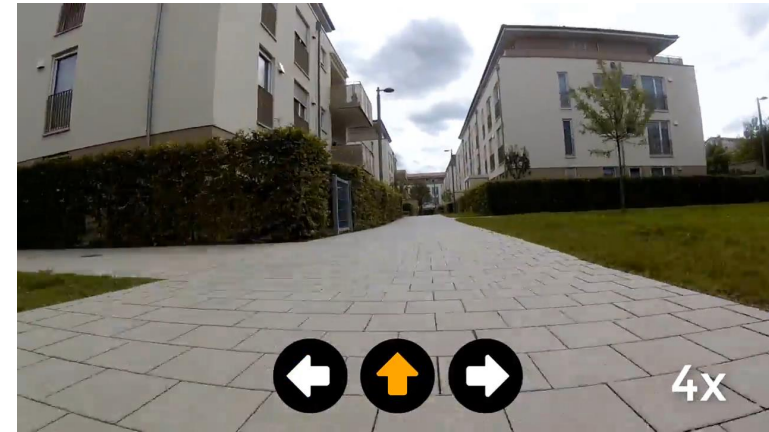
Gregory Kahn

Why should you care?

Safety



Flexibility



Outline

Topics

- Safety
- Flexibility

Algorithms

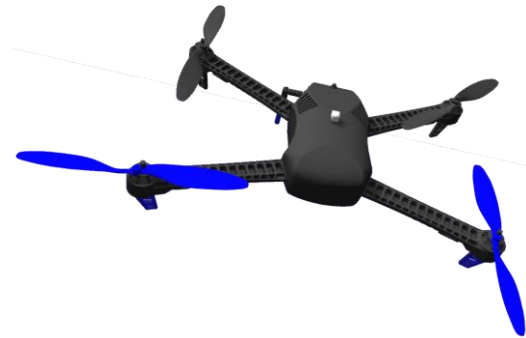
- Imitation learning
- Model-free
- Model-based

$2 * 3 = 6$ papers we'll cover

By no means the best / only papers on these topics

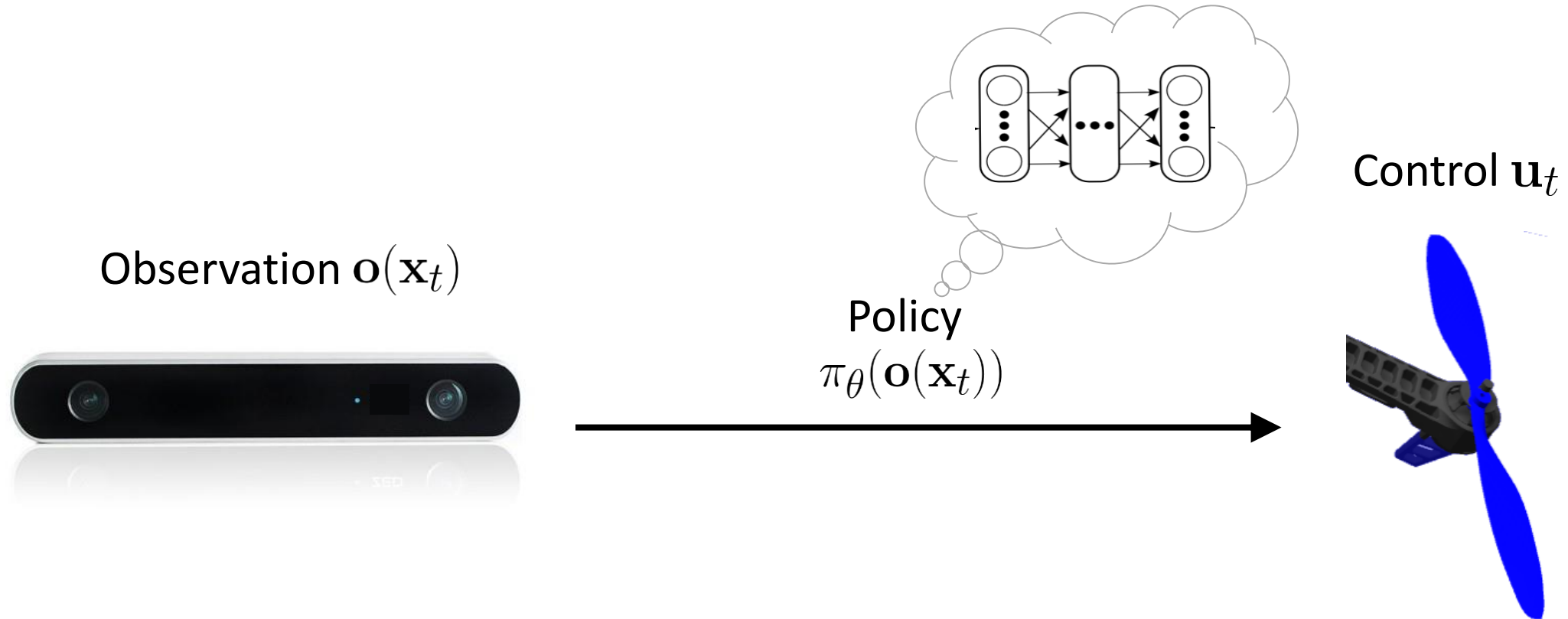
PLATO: Policy Learning using Adaptive Trajectory Optimization

Gregory Kahn¹, Tianhao Zhang¹, Sergey Levine¹, Pieter Abbeel^{1,2,3}



Goal

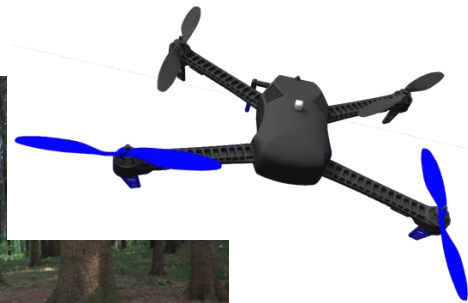
Learn control policy that **maps observations to controls**



Assumption

- Able to generate good trajectories using an expert policy π^*

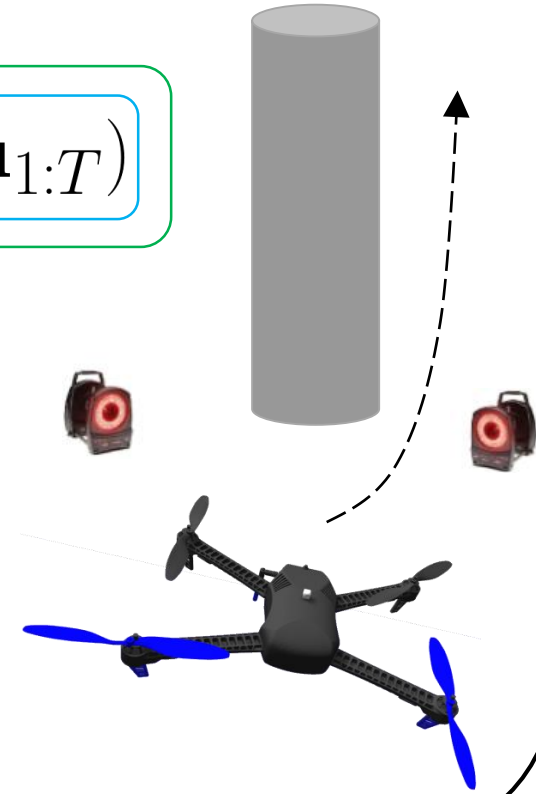
Human expert



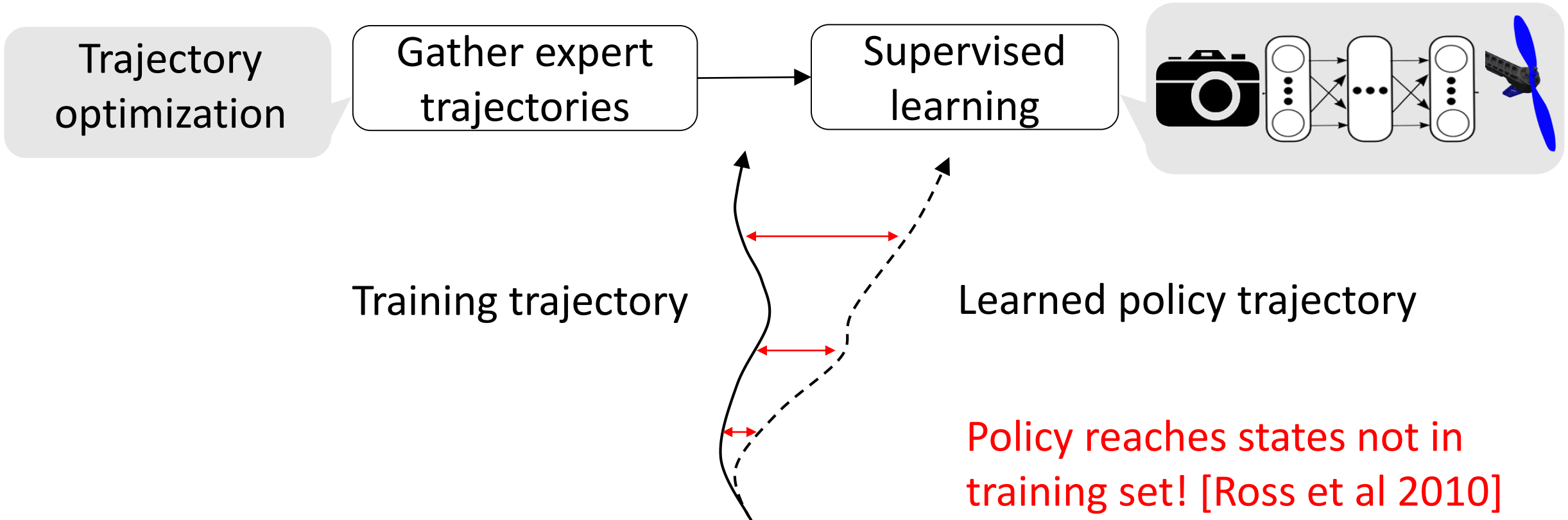
Trajectory optimization

$$c(\mathbf{x}_0, \mathbf{u}_{1:T})$$

- cost function
- optimization
- full state information
- only during training



Supervised Learning



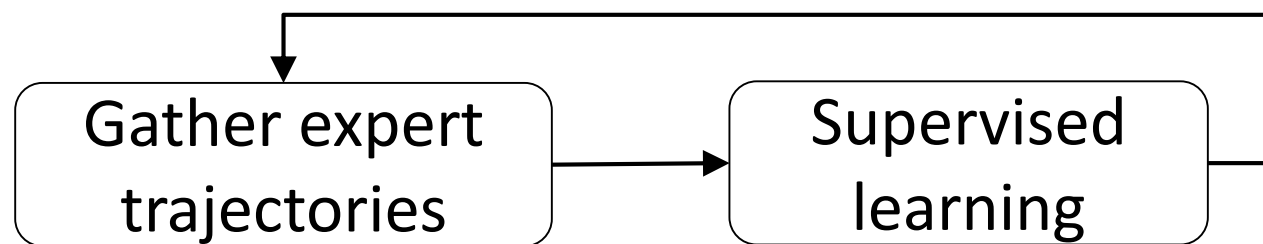
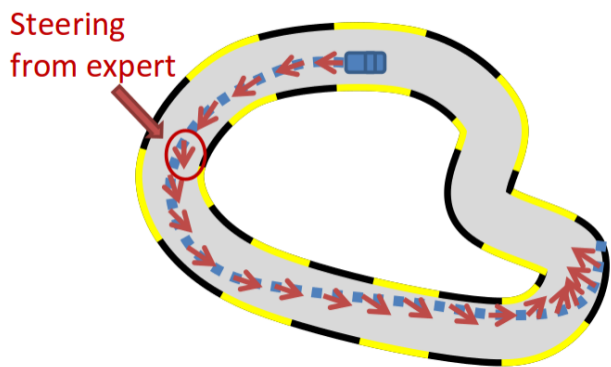
- Problem: training and test distributions differ

Dataset Aggregation (Dagger)

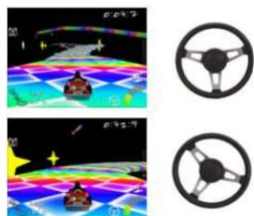
[Ross et al 2011]

- Problem: training and test distributions differ
- Solution: execute policy during training

$$\pi_{\text{mix}} \leftarrow \beta \pi^* + (1 - \beta) \pi_{\theta}$$

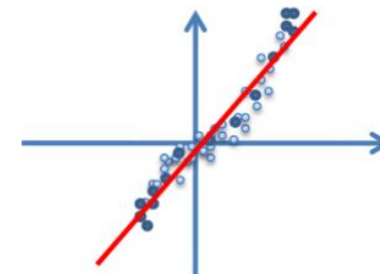
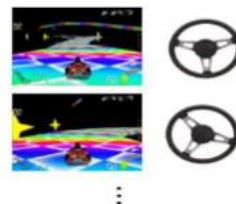


New Data



+

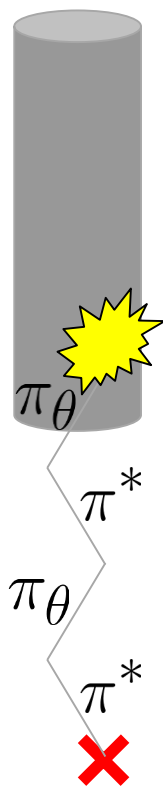
All previous data



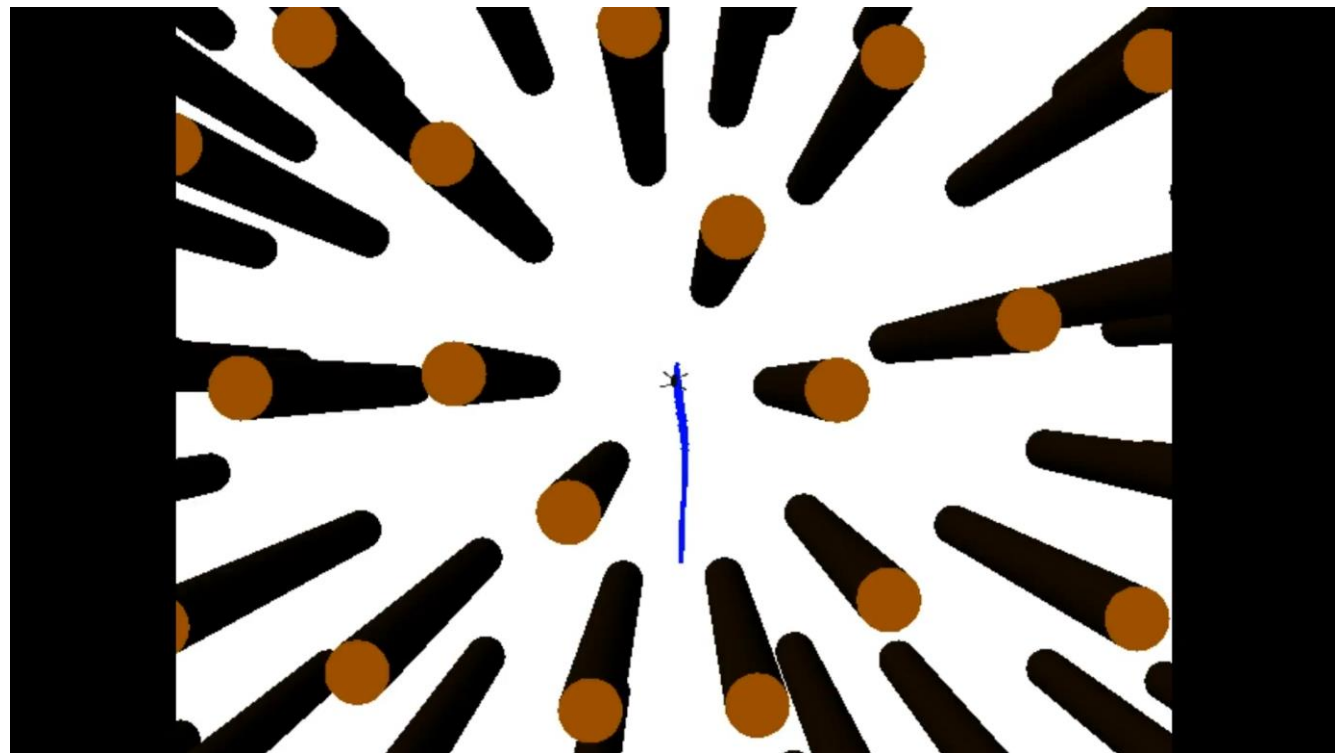
Safety during training

- DAgger mixes the actions

$$\mathbf{u}_{\text{mix}} \sim \begin{cases} \pi^*(\mathbf{u}|\mathbf{x}_t) & \text{prob. } \beta \\ \pi_\theta(\mathbf{u}|\mathbf{x}_t) & \text{prob. } (1 - \beta) \end{cases}$$



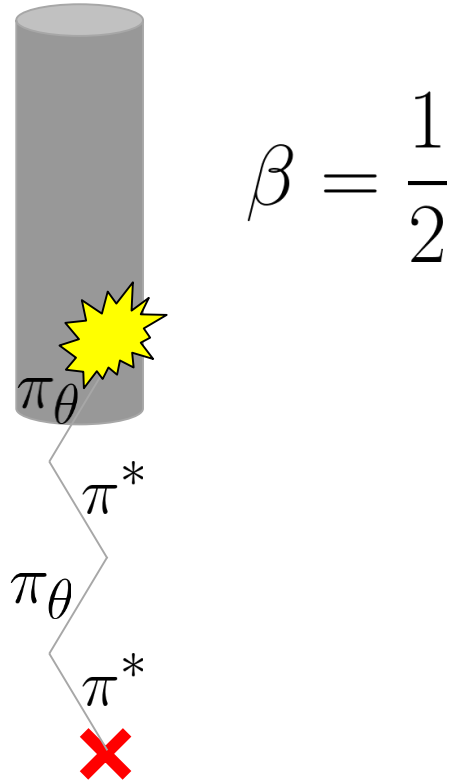
$$\beta = \frac{1}{2}$$



Policy Learning using Adaptive Trajectory Optimization (PLATO)

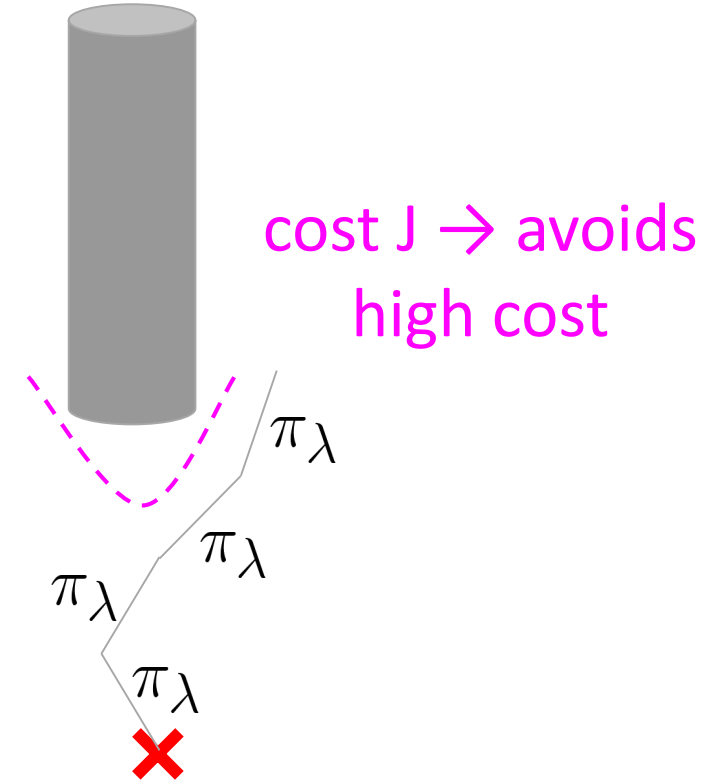
- DAgger mixes the actions

$$\mathbf{u}_{\text{mix}} \sim \begin{cases} \pi^*(\mathbf{u}|\mathbf{x}_t) & \text{prob. } \beta \\ \pi_\theta(\mathbf{u}|\mathbf{x}_t) & \text{prob. } (1 - \beta) \end{cases}$$



- PLATO mixes the objectives

$$\pi_\lambda \leftarrow \arg \min_{\pi} J(\pi) + \lambda D_{\text{KL}}(\pi || \pi_\theta)$$
$$\mathbf{u}_\lambda \sim \pi_\lambda(\mathbf{u}|\mathbf{x}_t)$$



Algorithm comparisons

approach	sampling policy	safe	similar training and test distributions
supervised learning	π^*	✓	✗
DAgger	π_{mix}	✗	✓
PLATO	π_λ	✓	✓

$$\pi_{\text{mix}} \leftarrow \beta \pi^* + (1 - \beta) \pi_\theta$$

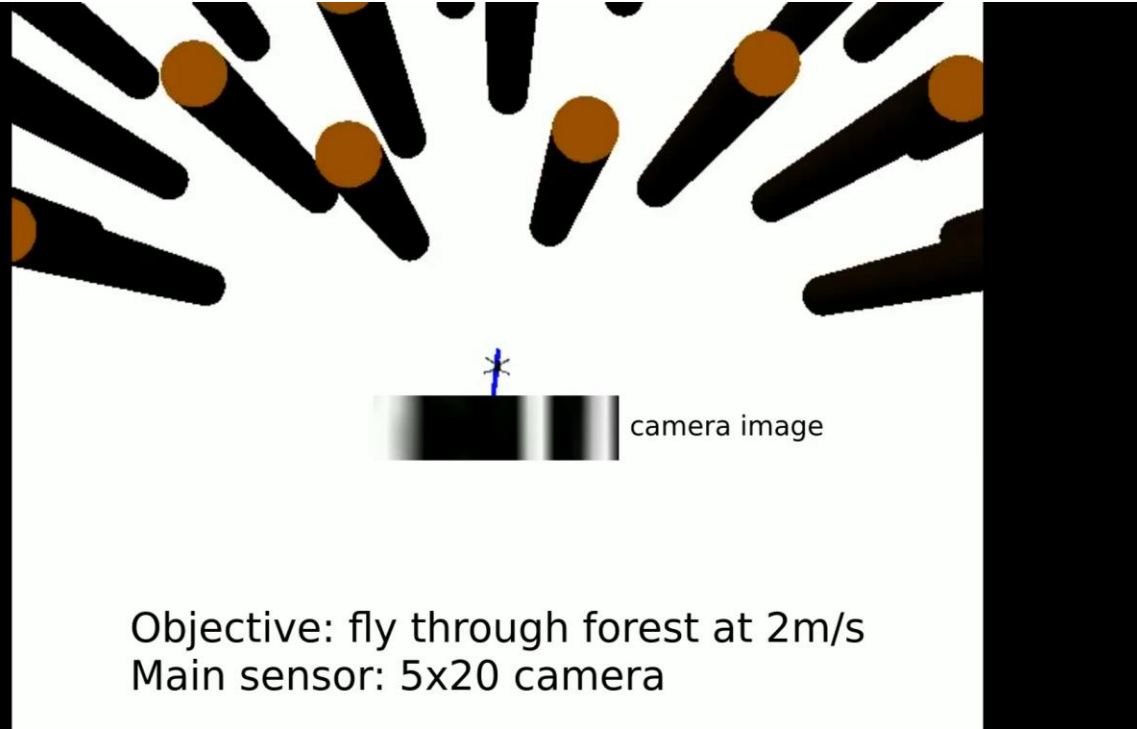
$$\pi_\lambda \leftarrow \arg \min_{\pi} J(\pi) + \lambda D_{\text{KL}}(\pi || \pi_\theta)$$

Experiments: final neural network policies

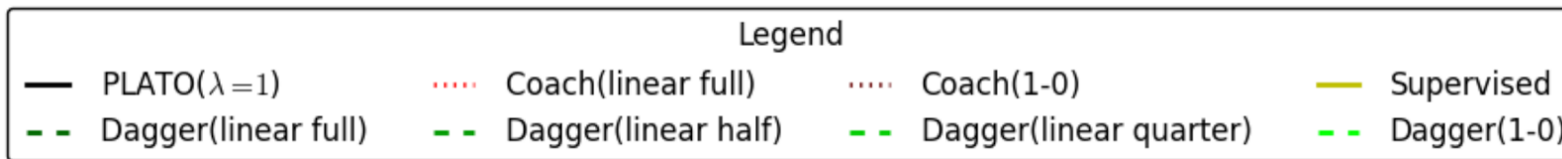
Canyon



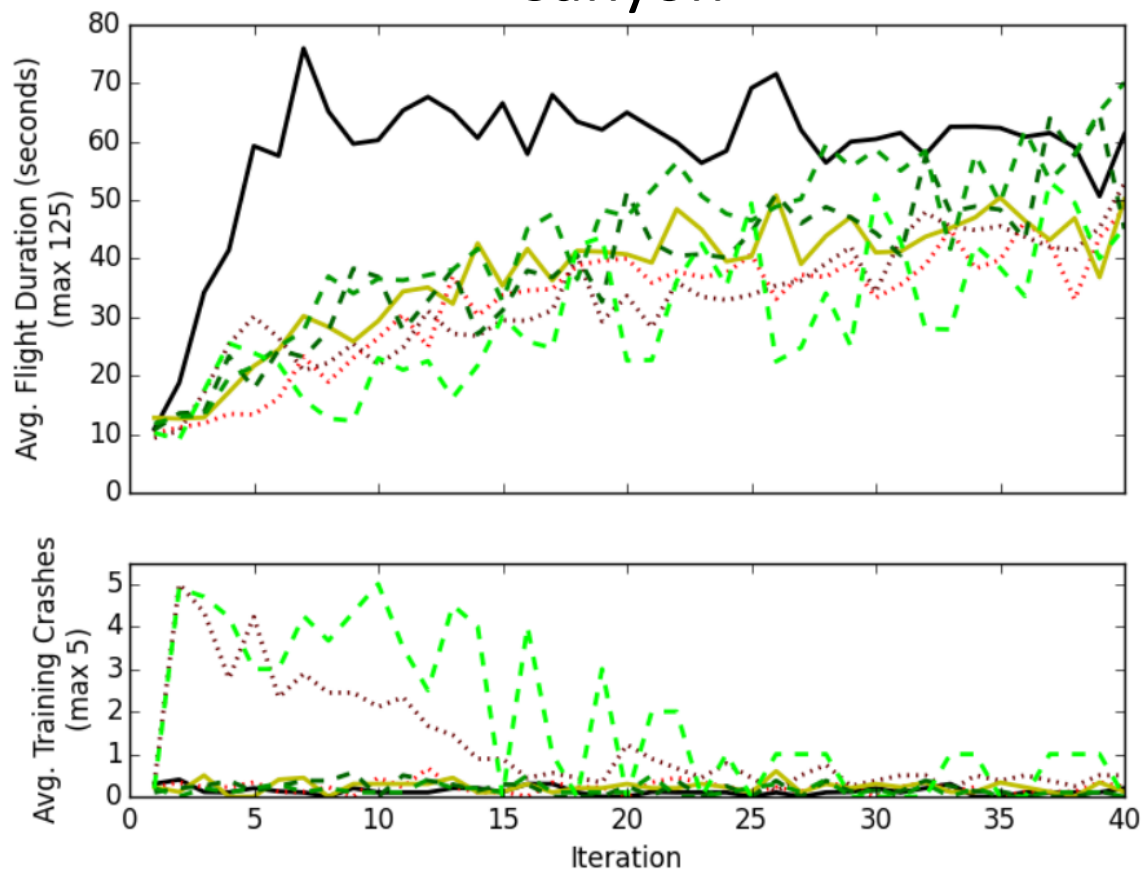
Forest



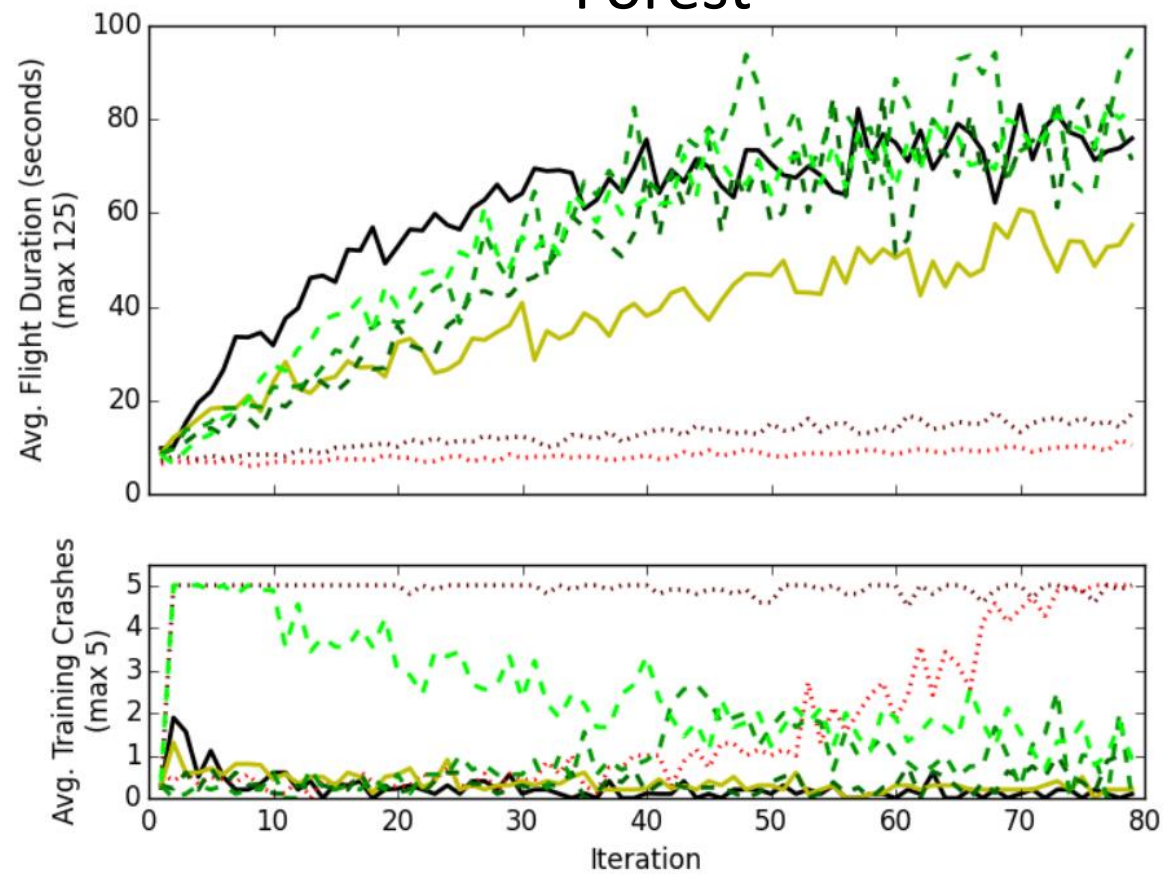
Experiments: metrics



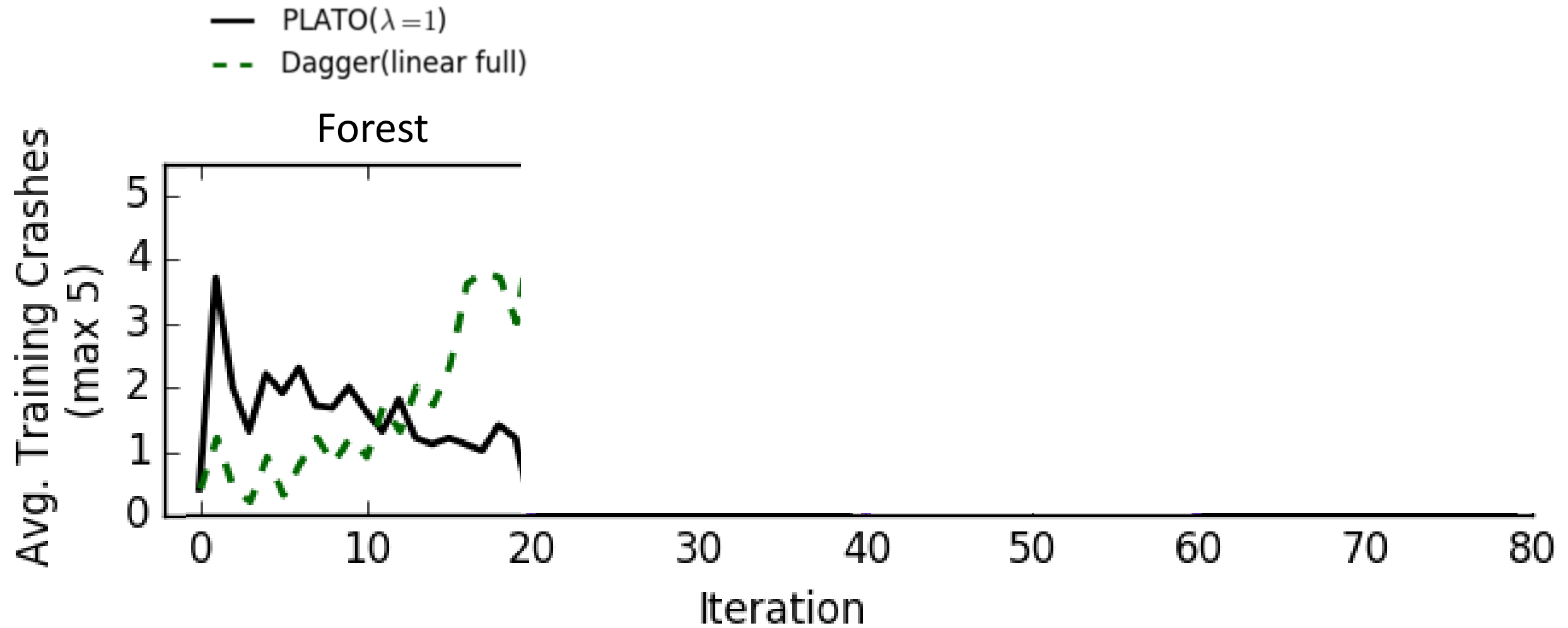
Canyon



Forest



Experiments: metrics

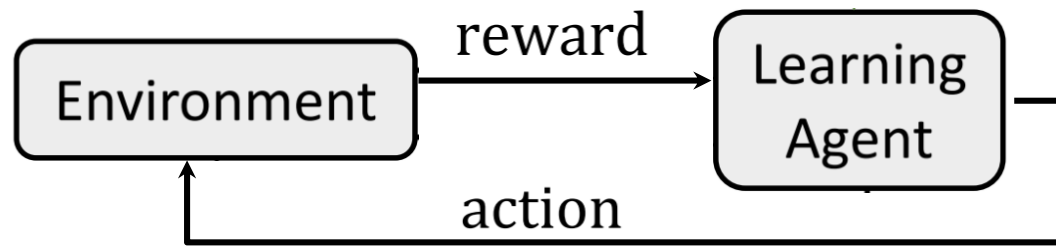


Safe Reinforcement Learning via Shielding

Mohammed Alshiekh¹, Roderick Bloem², Rüdiger Ehlers³, Bettina Könighofer², Scott Niekum¹, Ufuk Topcu¹

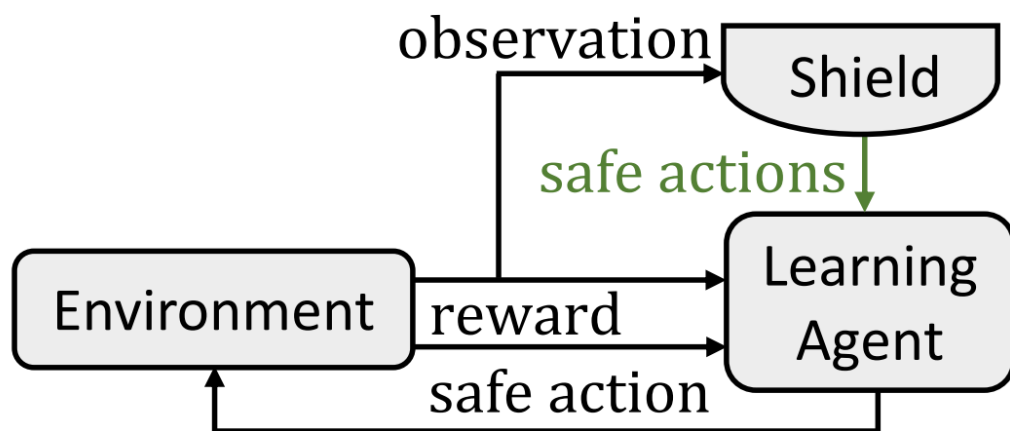
Goal

NOT SAFE



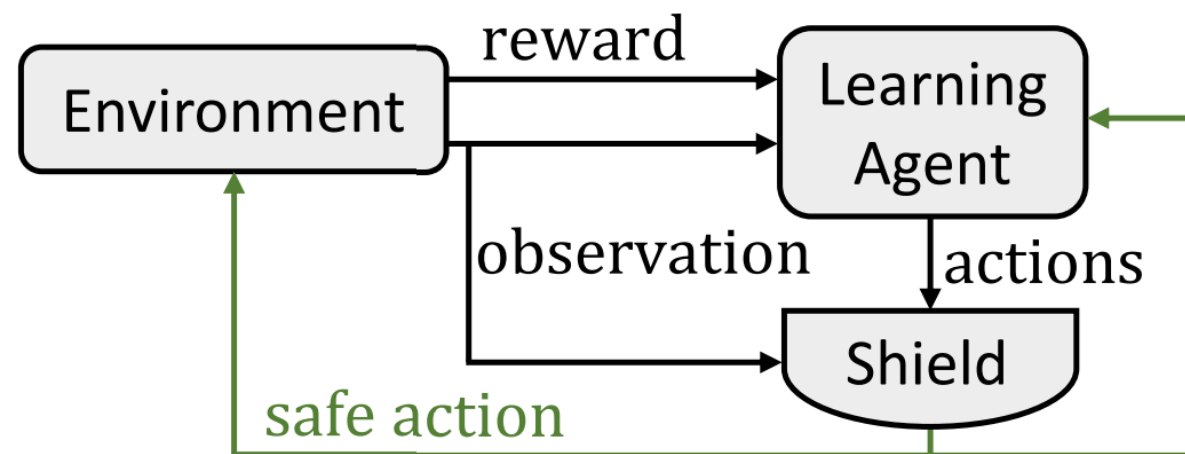
Shielding

Pre-emptive shielding



Like learning in a transformed MDP

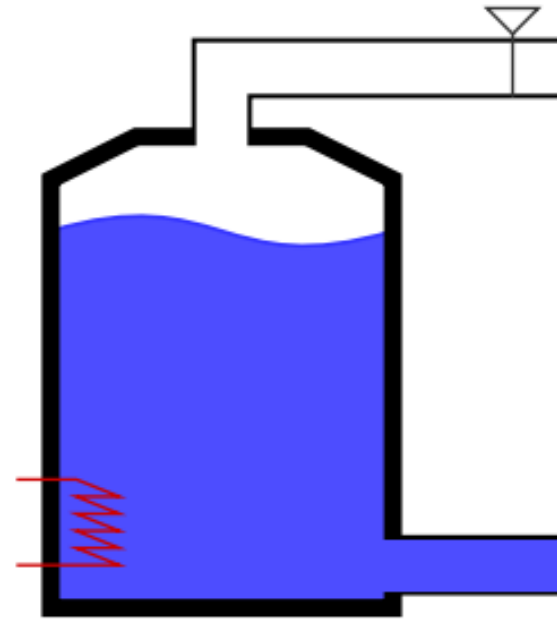
Post-posed shielding



Shield can be used at test time

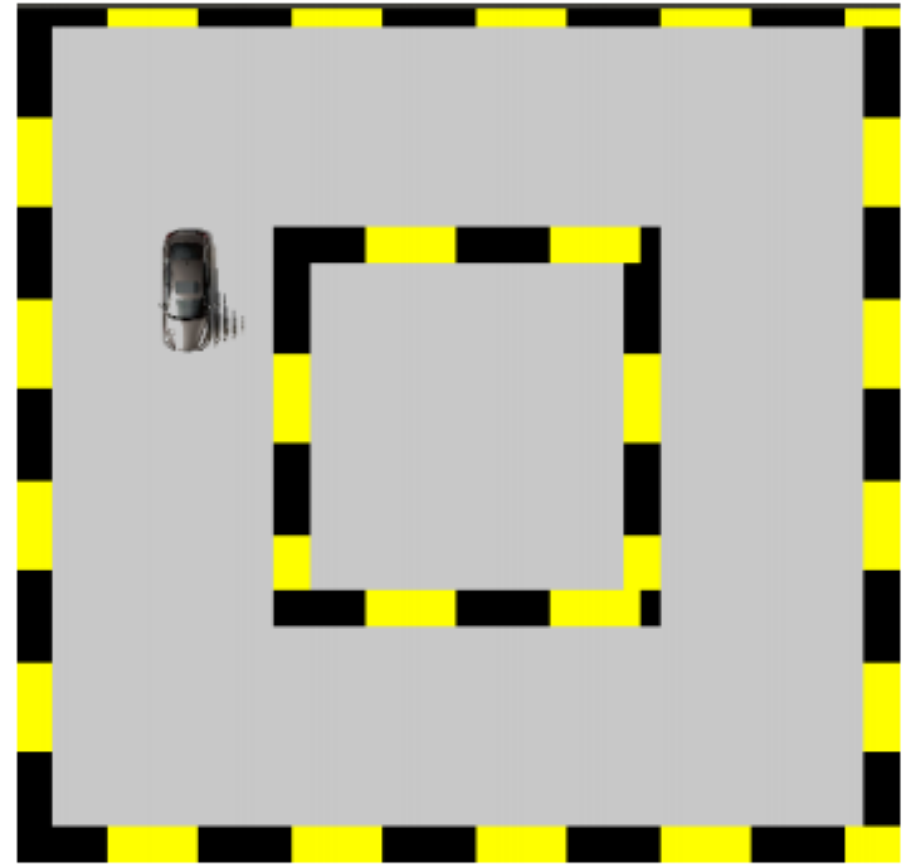
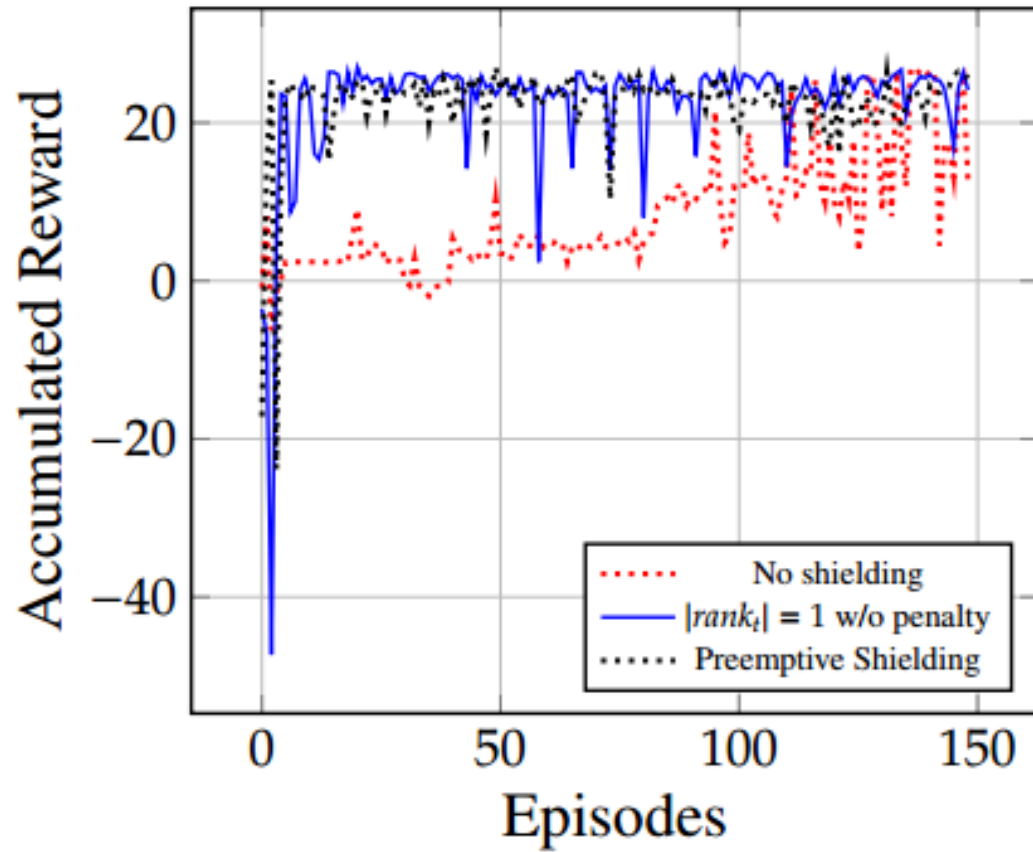
How to shield: linear temporal logic

- Encode safety with temporal logic
- Assumption: Known approximate/conservative transition dynamics

$$\begin{aligned} & G(\text{level} > 0) \\ \wedge & G(\text{level} < 100) \\ \wedge & G((\text{open} \wedge X\text{close}) \rightarrow XX\text{close} \wedge XXX\text{close}) \\ \wedge & G((\text{close} \wedge X\text{open}) \rightarrow XX\text{open} \wedge XXX\text{open}) \end{aligned}$$


Experiments

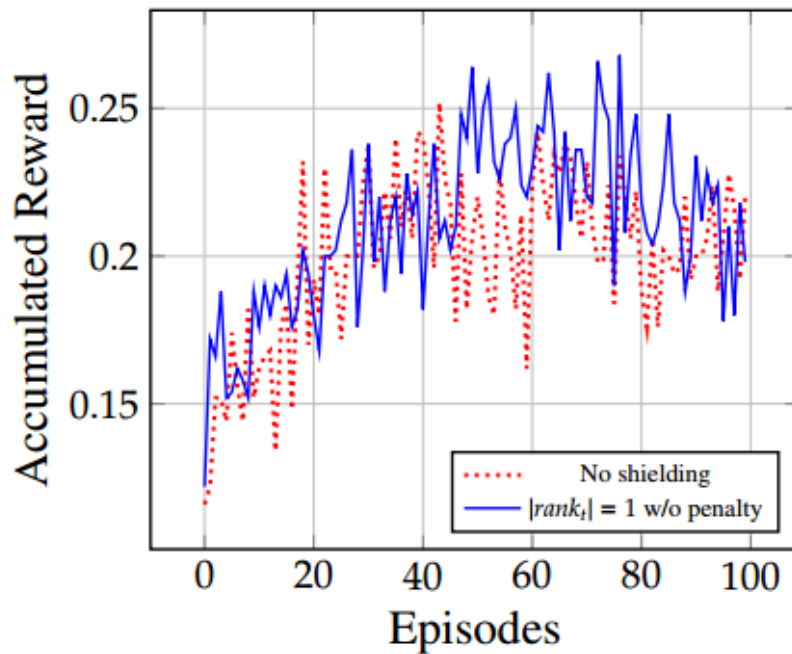
Safety criteria
- Don't crash



Experiments

Safety criteria

- Don't run out of oxygen
- If enough oxygen, don't surface w/o divers

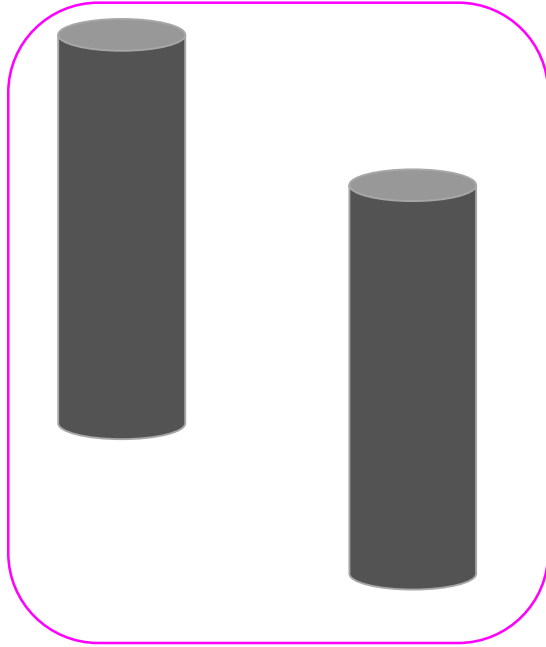


Uncertainty-Aware Reinforcement Learning for Collision Avoidance

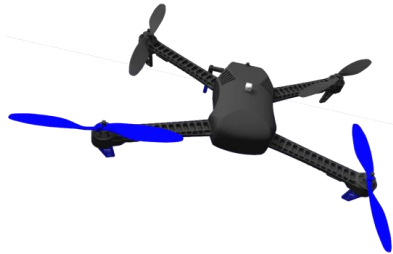
Gregory Kahn*, Adam Villaflor*, Vitchyr Pong*, Pieter Abbeel*[†], Sergey Levine*

Goal

unknown environment

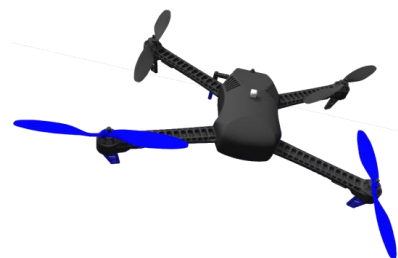
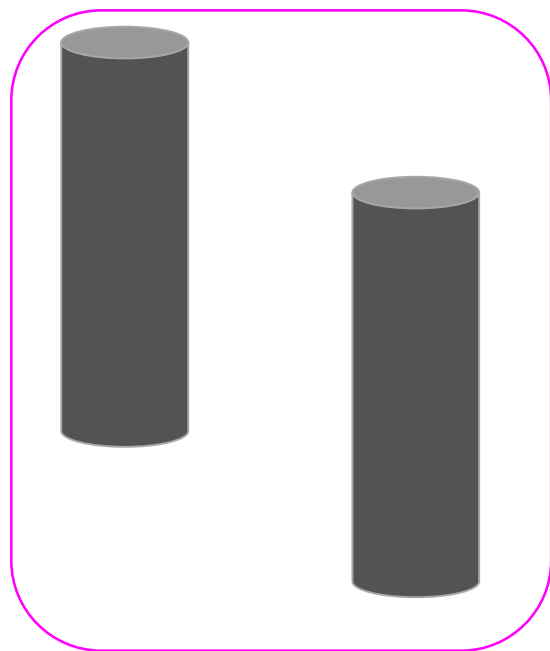


How to do reinforcement learning
without destroying the robot during training
using only onboard images



Approach

unknown environment



$$c(\tau) = c_{\text{TASK}}(\tau) + c_{\text{COLL}}(\tau)$$

learn a collision prediction model

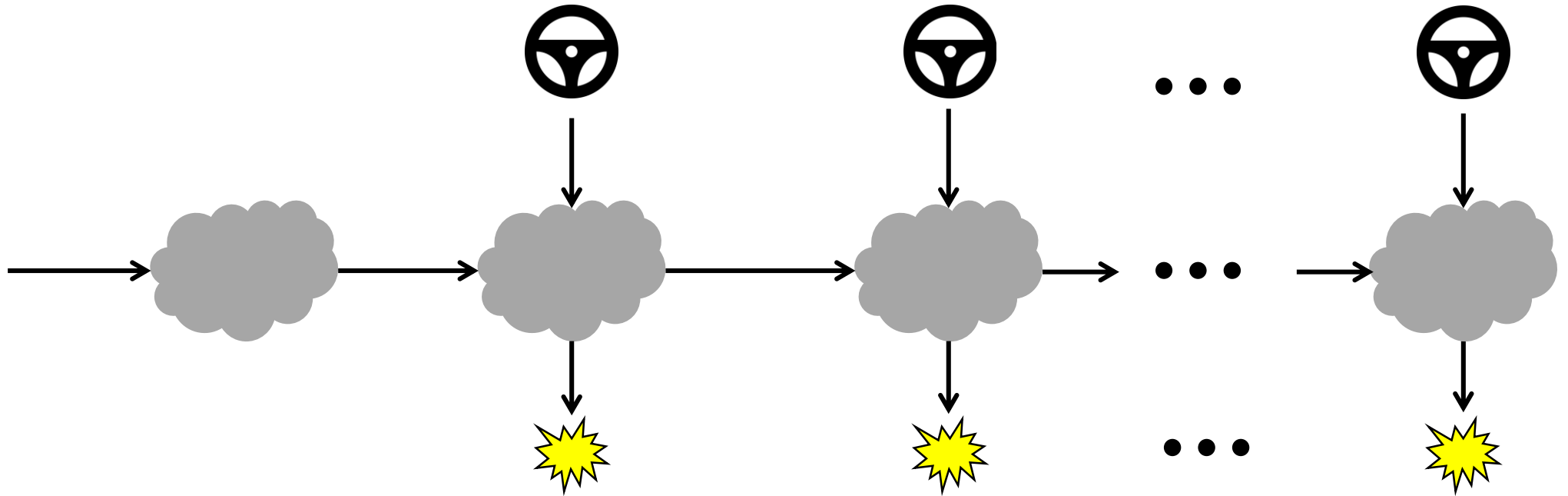
$$p(c_{t+H} | \mathbf{o}_t, \mathbf{u}_t, \dots, \mathbf{u}_{t+H})$$

raw image

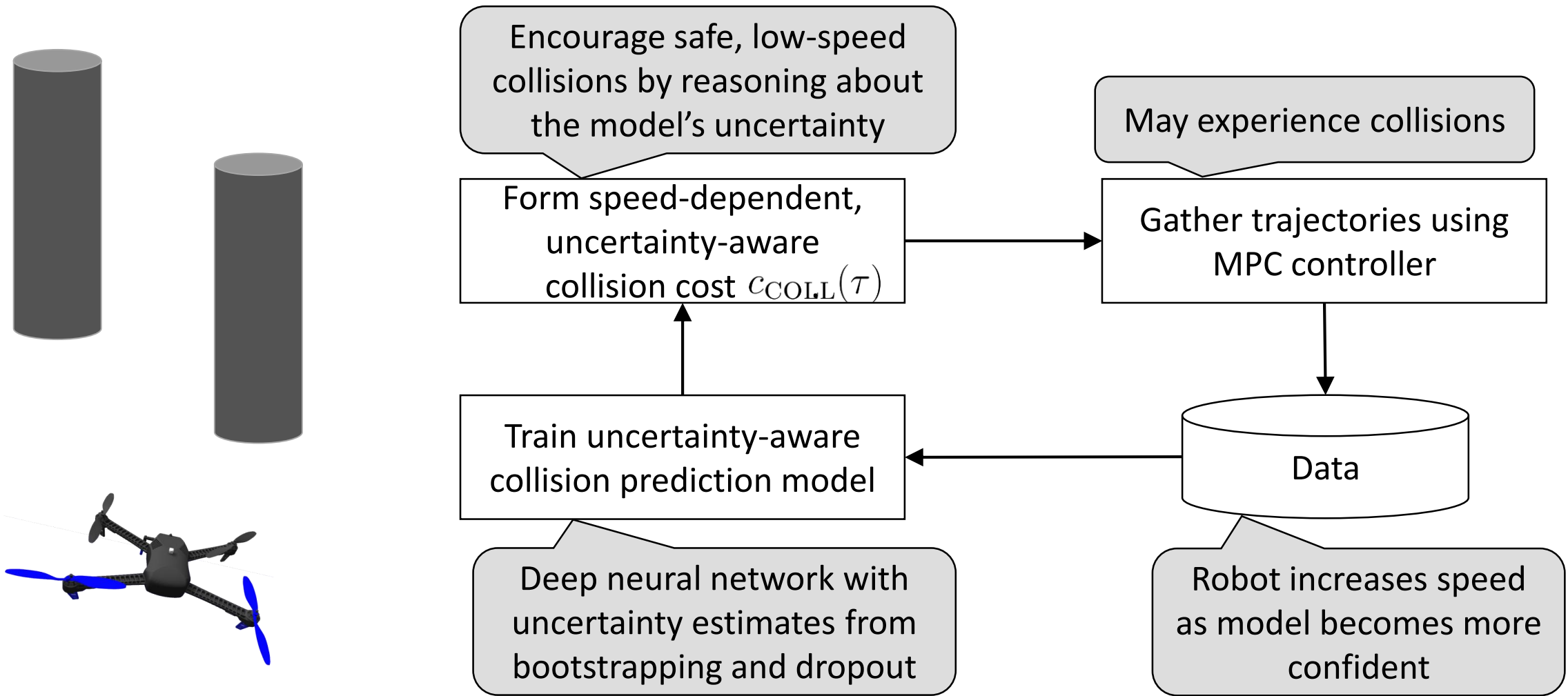
command velocities

neural network

Collision prediction model



Model-based RL using collision prediction model



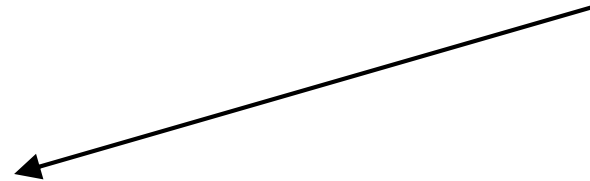
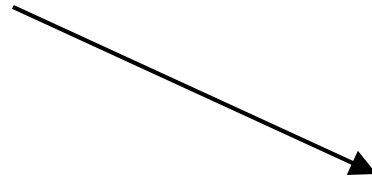
Collision cost

$$c_{\text{COLL}}(\tau) \propto \text{SPEED} \cdot \left(\text{E}[p(c_{t+H} | \tau)] + \sqrt{\text{Var}[p(c_{t+H} | \tau)]} \right)$$

high speed

predict collision

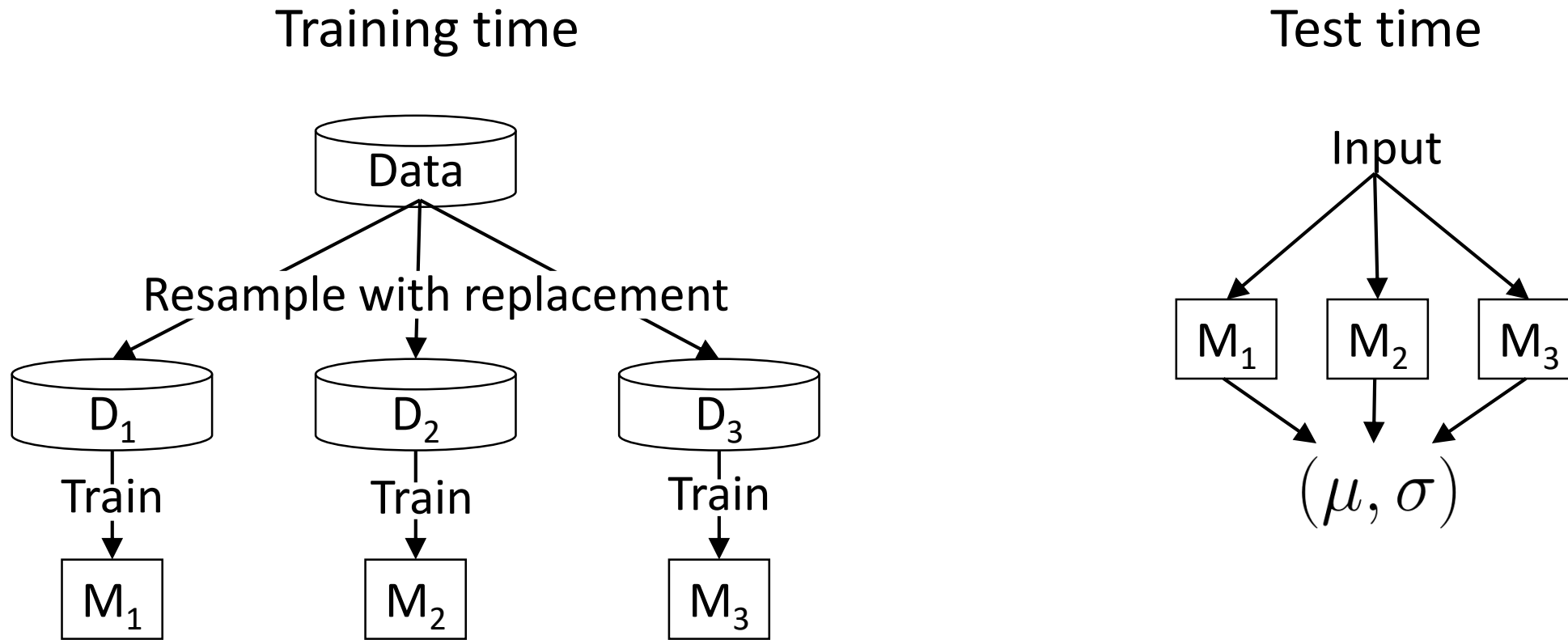
large uncertainty



large cost

Estimating neural network output uncertainty

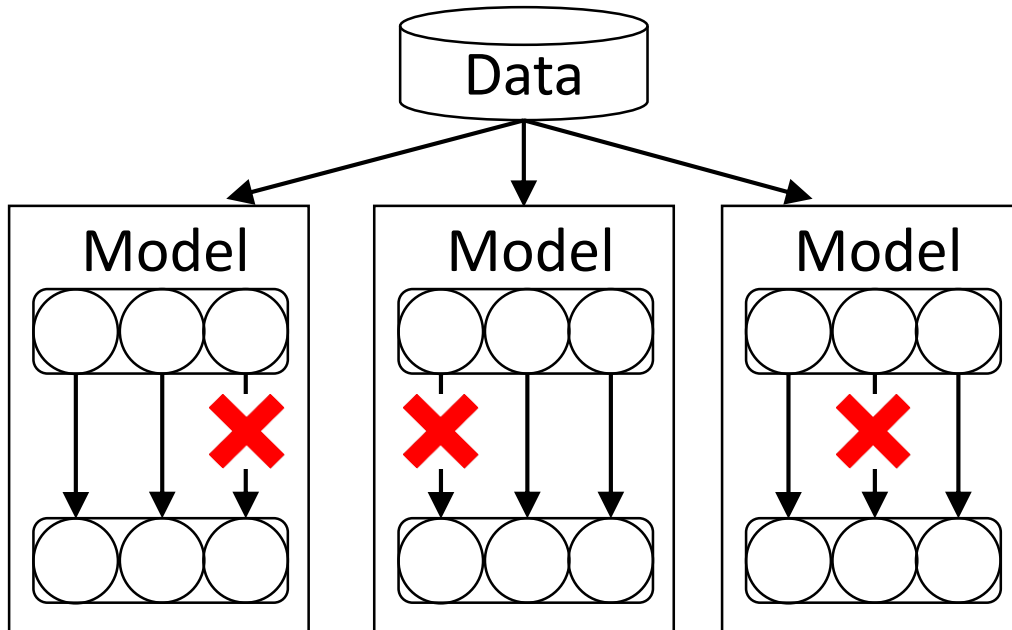
Bootstrapping



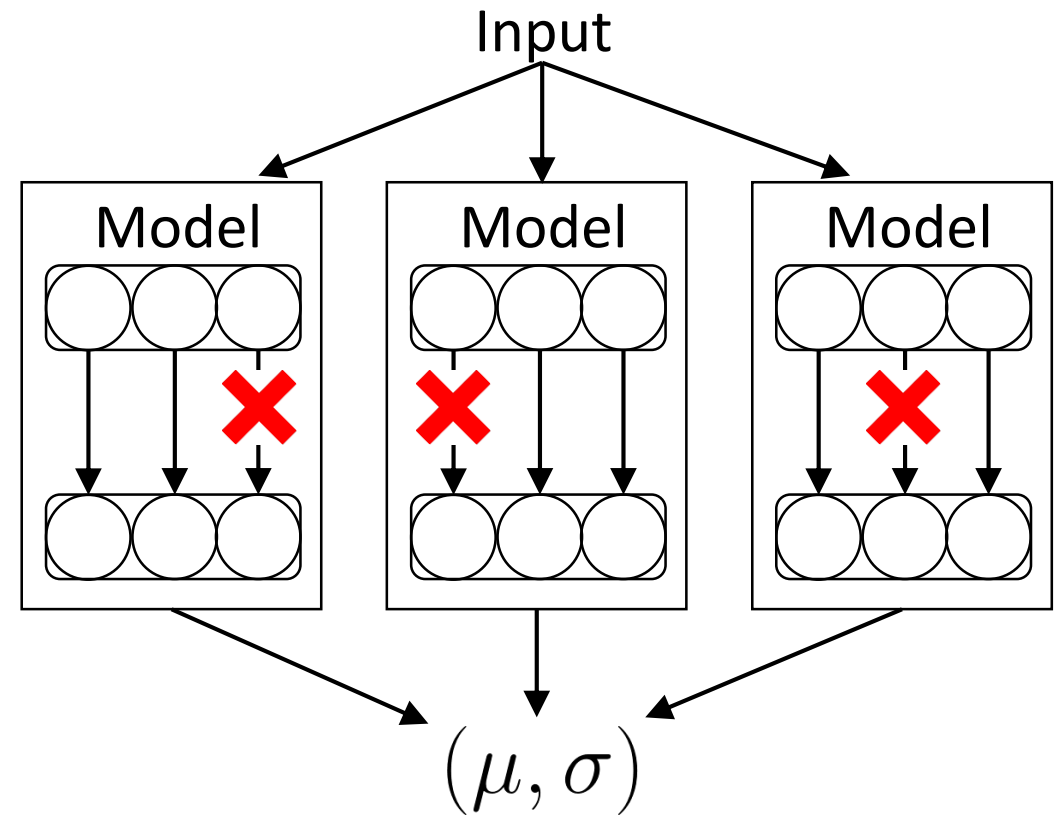
Estimating neural network output uncertainty

Dropout

Training time

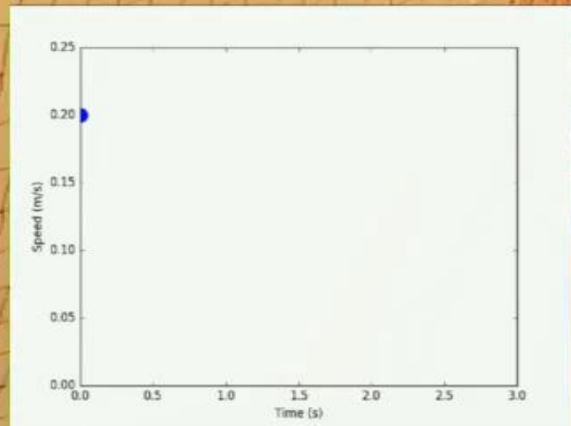


Test time



Preliminary real-world experiments

Not accounting for uncertainty
(higher-speed collisions)



Preliminary real-world experiments

accounting for uncertainty
(lower-speed collisions)



Preliminary real-world experiments

successful flight past obstacle



Safety takeaways

- Tradeoff between safety and exploration
- Safety guarantees require expert oversight or known environment + dynamics
- Uncertainty can play a key role

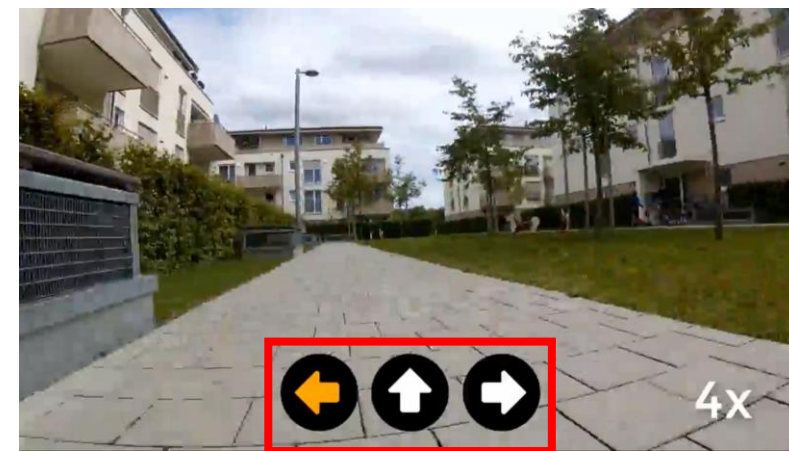
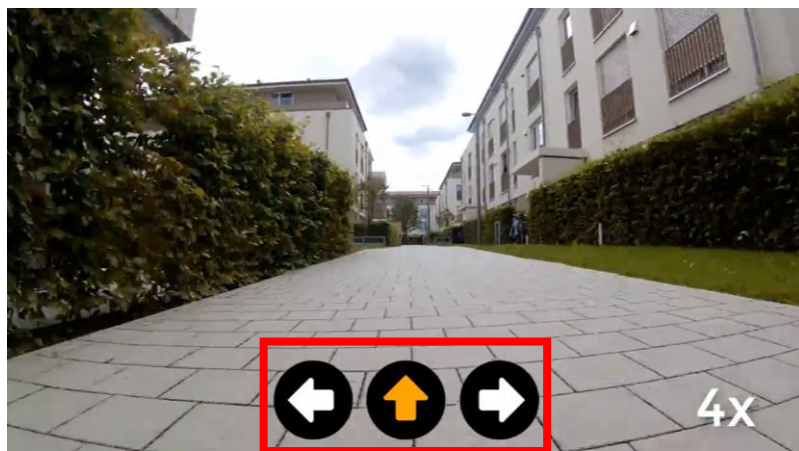
End-to-end Driving via Conditional Imitation Learning

Felipe Codevilla^{1,2} Matthias Müller^{1,3} Antonio López² Vladlen Koltun¹ Alexey Dosovitskiy¹

Goal

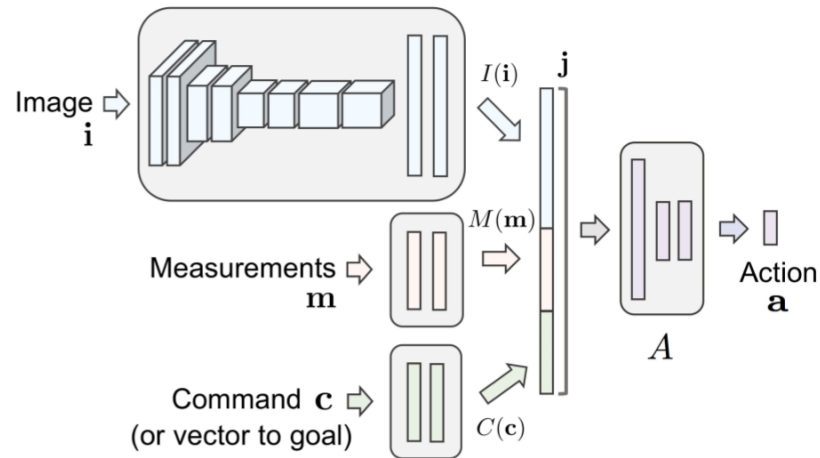
$$\min_{\theta} \|\pi_{\theta}(\mathbf{a}|\mathbf{s}) - \pi^*(\mathbf{a}|\mathbf{s})\| \quad \longrightarrow \quad \min_{\theta} \|\pi_{\theta}(\mathbf{a}|\mathbf{s}, \underline{\mathbf{c}}) - \pi^*(\mathbf{a}|\mathbf{s}, \underline{\mathbf{c}})\|$$

User-specified command

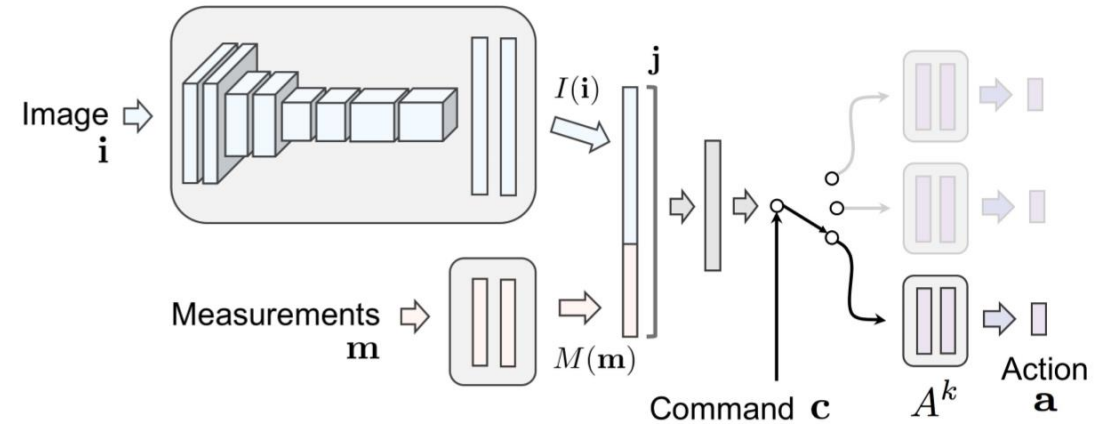


Approach

Option A: Input command



Option B: Branch using command



+ empirically better

- only works for discrete commands

Approach

Important details

- Data augmentation
 - Contrast
 - Brightness
 - Tone
 - Gaussian blur
 - Salt-and-pepper noise
 - Region dropout
- Adding noise to expert



4x



4x

Generalization

We evaluate how our model generalizes to previously unseen environments with very different appearance.

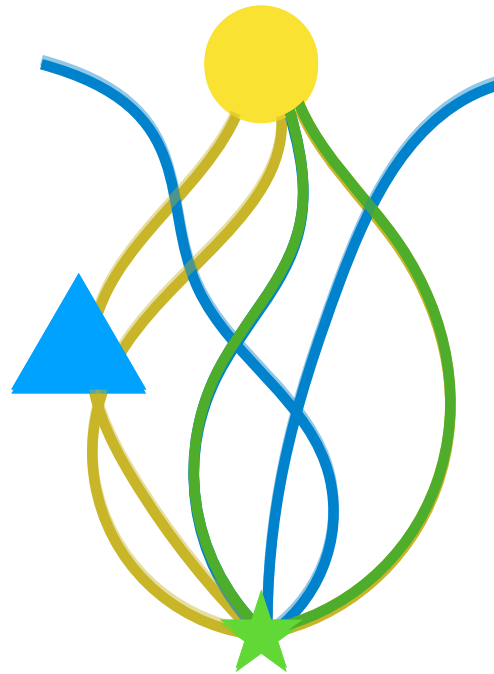
Composable Deep Reinforcement Learning for Robotic Manipulation

Tuomas Haarnoja¹, Vitchyr Pong¹, Aurick Zhou¹, Murtaza Dalal¹, Pieter Abbeel^{1,2}, Sergey Levine¹

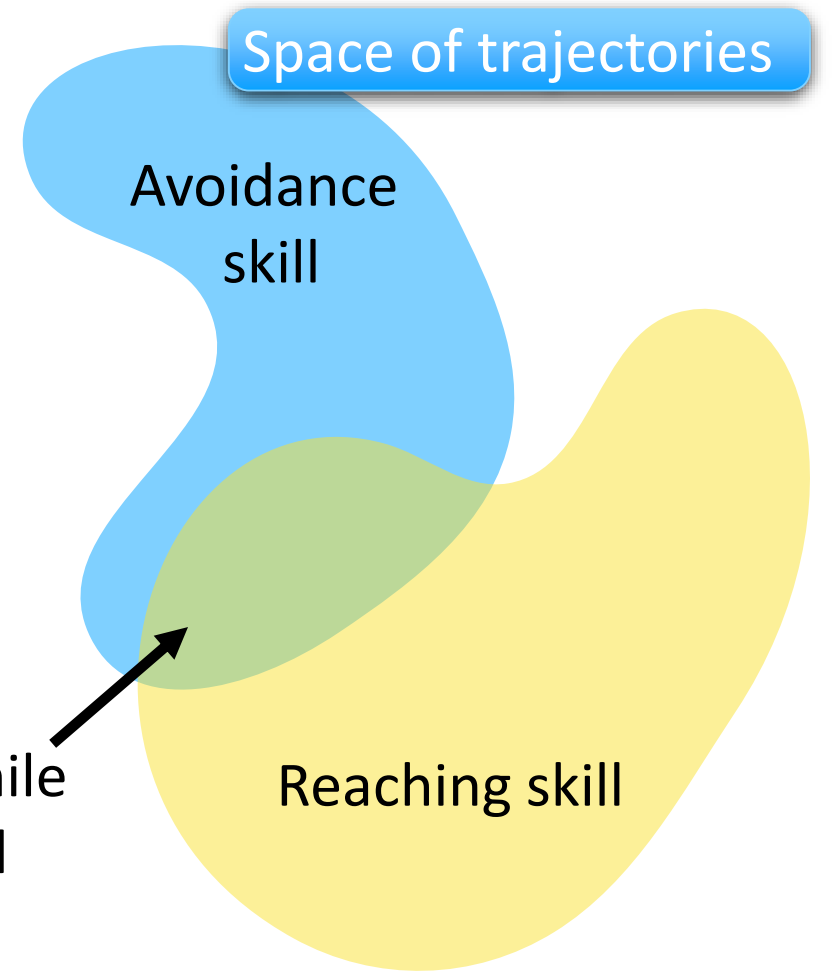
[slides adapted from Tuomas Haarnoja]

Goal

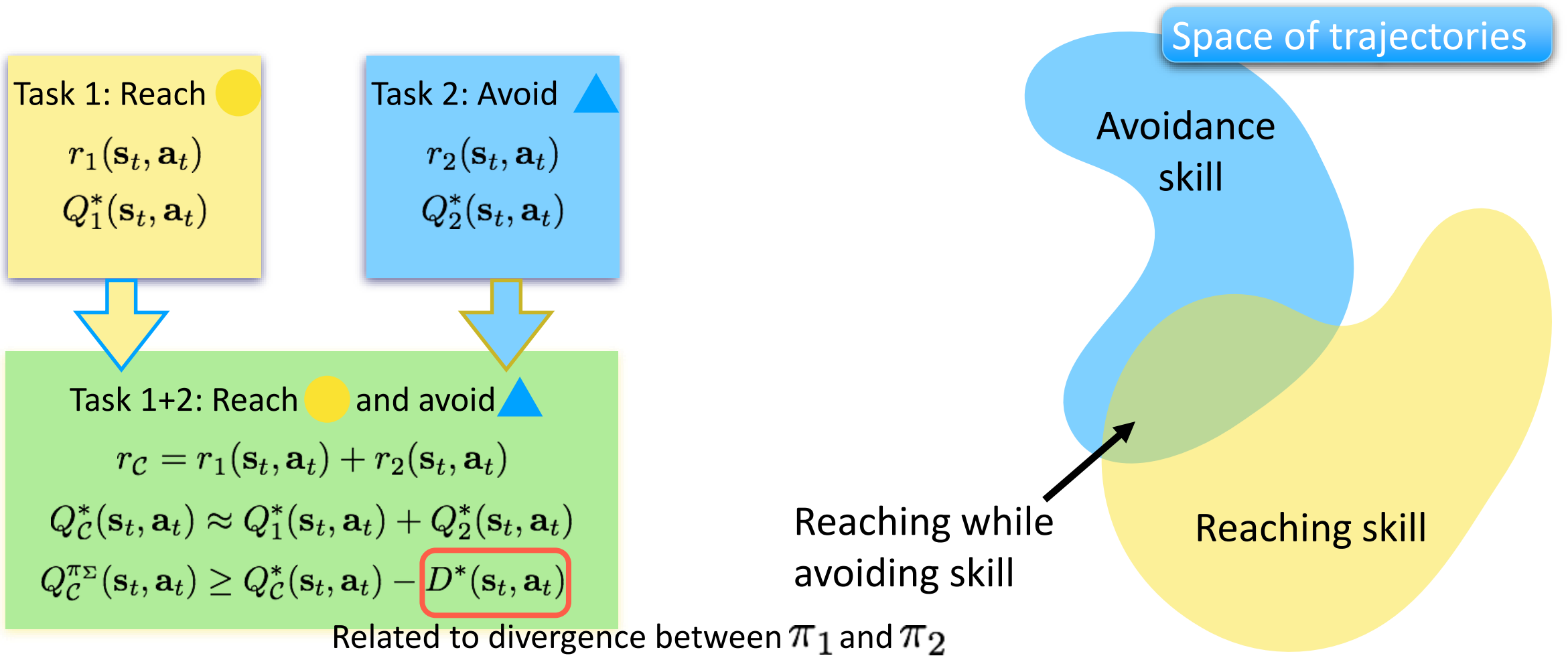
Task 1: Reach ●
Task 2: Avoid ▲



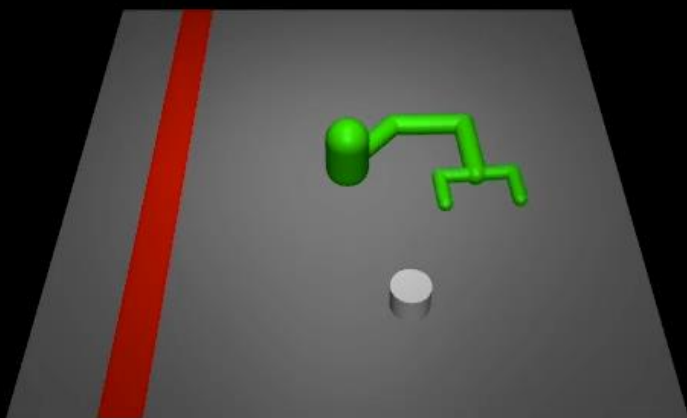
Reaching while
avoiding skill



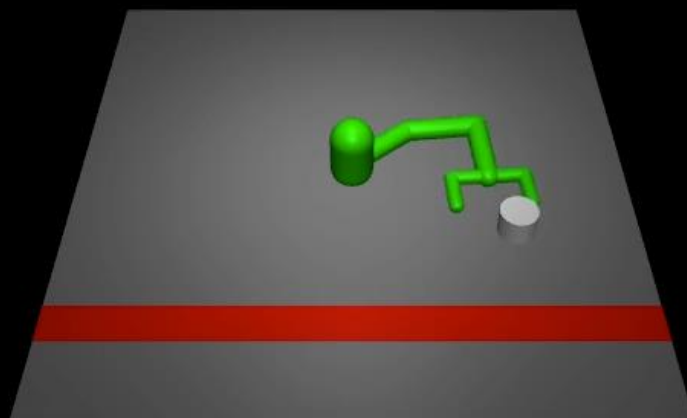
Policy Composition



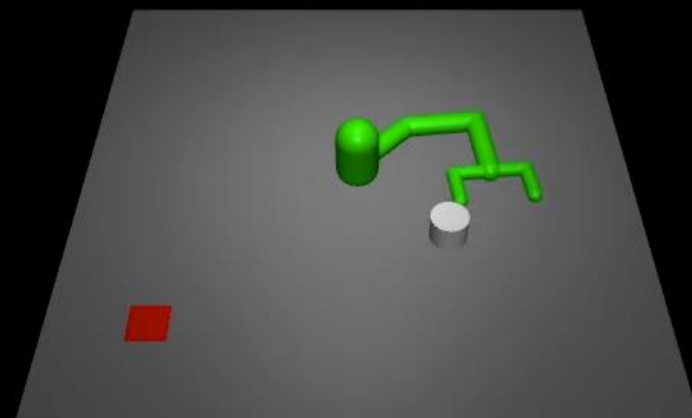
Reusability!



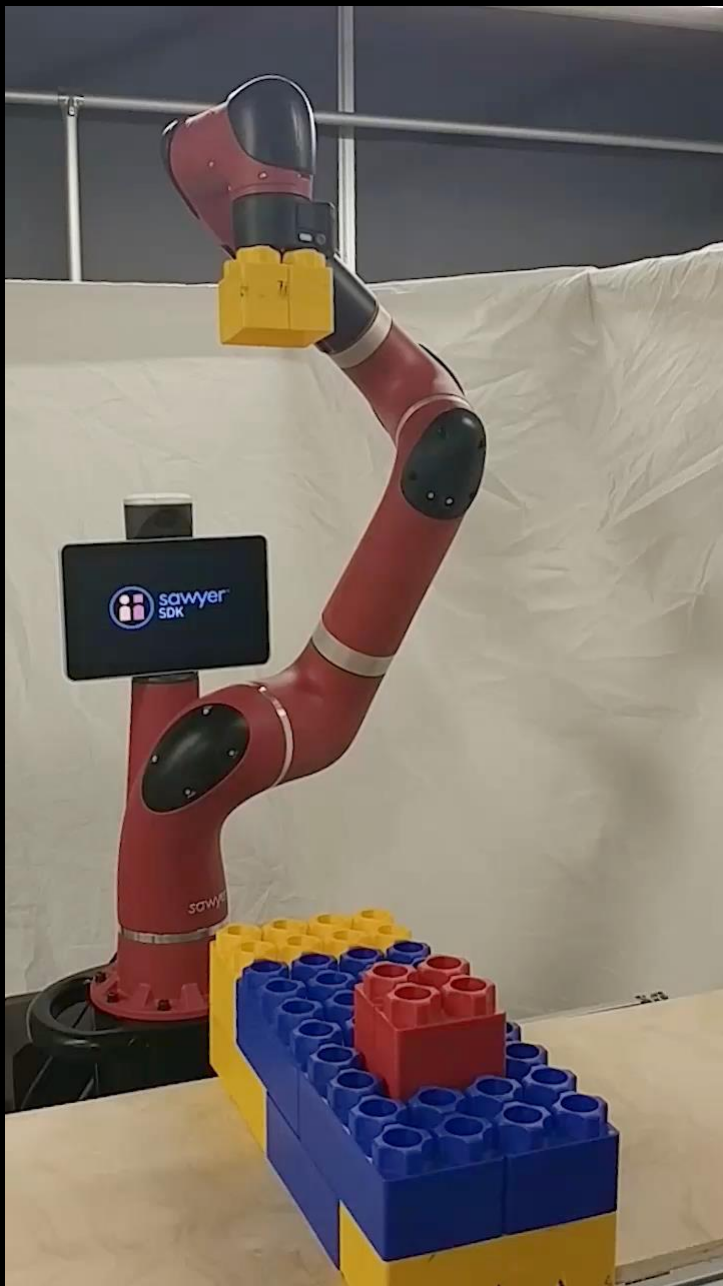
Task 1
 $Q_1^*(\mathbf{s}_t, \mathbf{a}_t)$



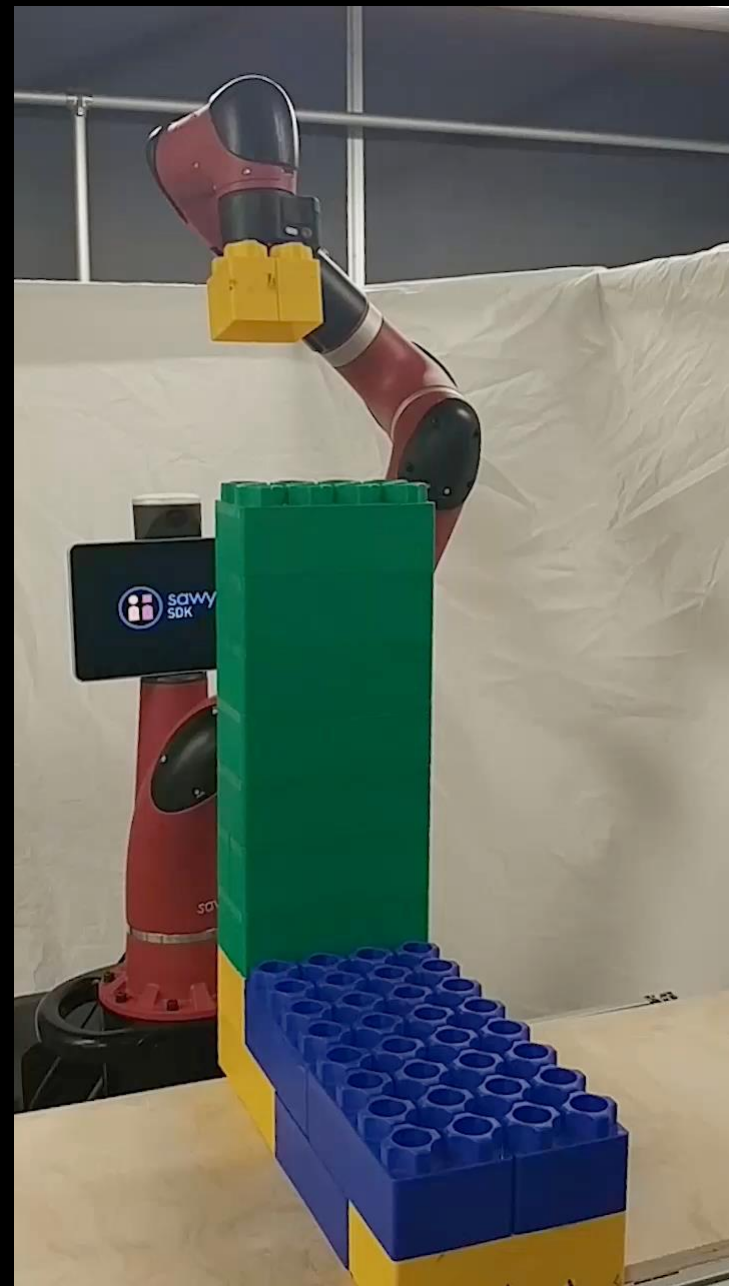
Task 2
 $Q_2^*(\mathbf{s}_t, \mathbf{a}_t)$



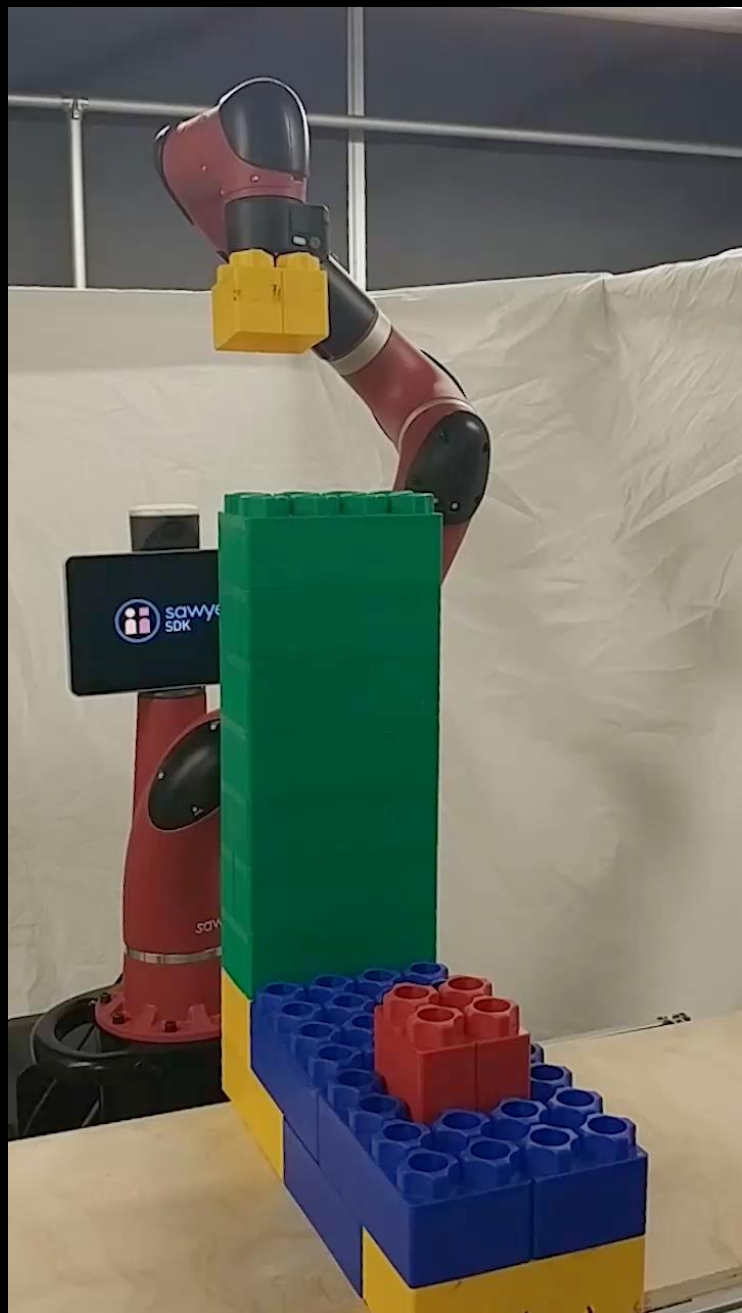
Task 1 + 2
 $Q_1^*(\mathbf{s}_t, \mathbf{a}_t) + Q_2^*(\mathbf{s}_t, \mathbf{a}_t)$



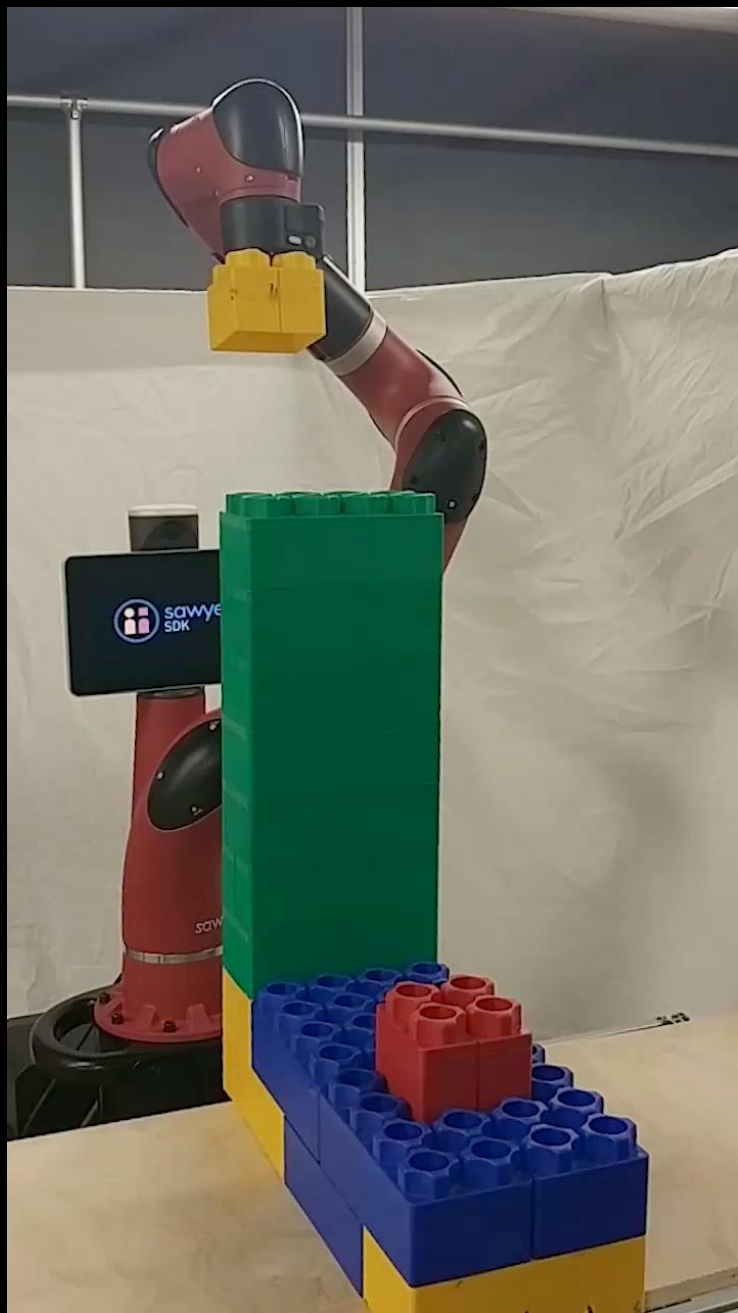
Stacking policy



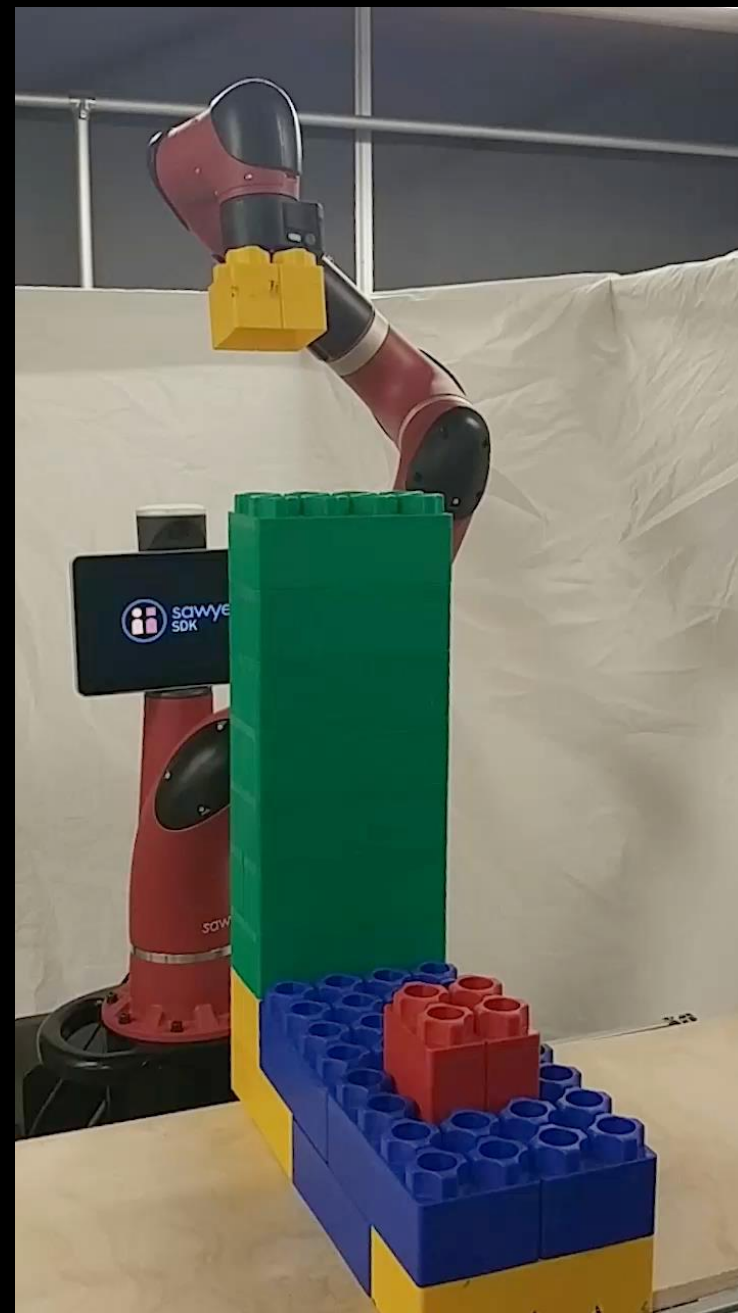
Avoidance policy



Stacking policy



Avoidance policy



Combined policy

Composable Action-Conditioned Predictors: Flexible Off-Policy Learning for Robot Navigation

Gregory Kahn*, Adam Villaflor*, Pieter Abbeel, Sergey Levine

Standard Reinforcement Learning

Train

Test

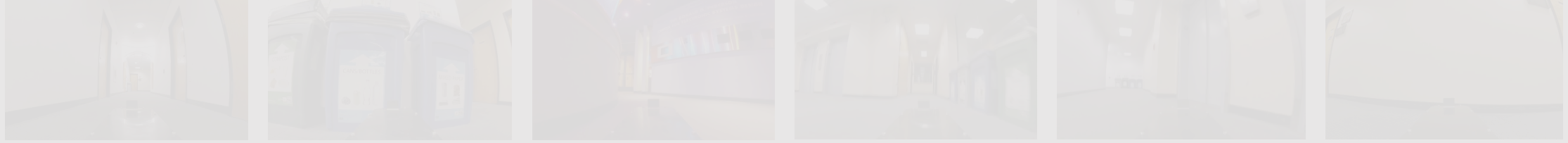


Data inefficient
Expert in the loop
Inflexible



CAPs Approach

Train



Data efficient
Detector in the loop
Flexible



Test



Detect

Predict

Control

Safety

Flexibility

Imitation learning

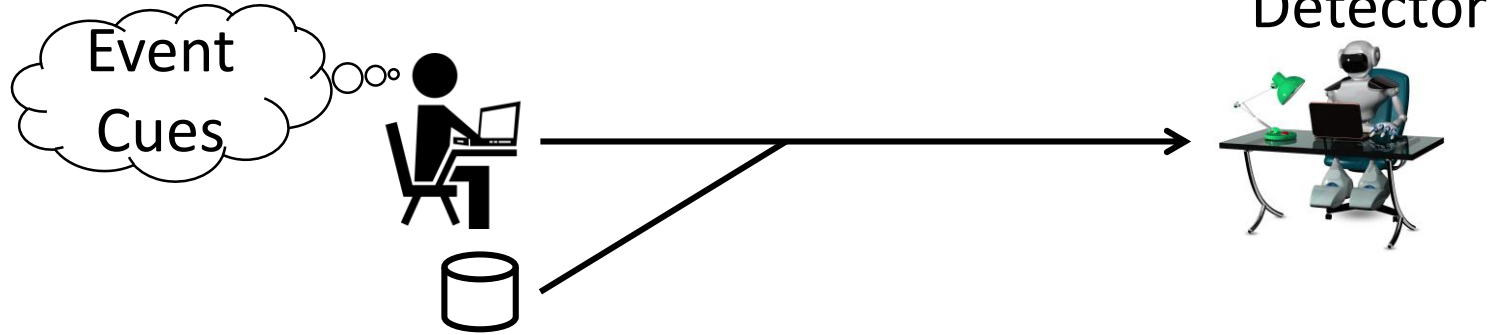
Model-free

Model-based

Detect

Predict

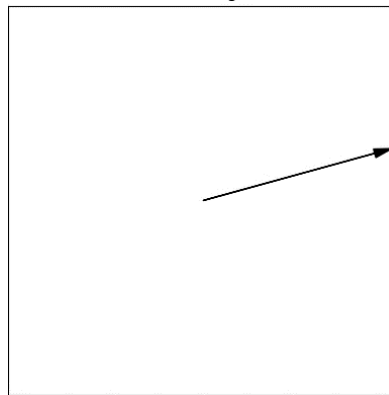
Control



$e_{\text{COLLISION}}$



e_{HEADING}



e_{DOOR}

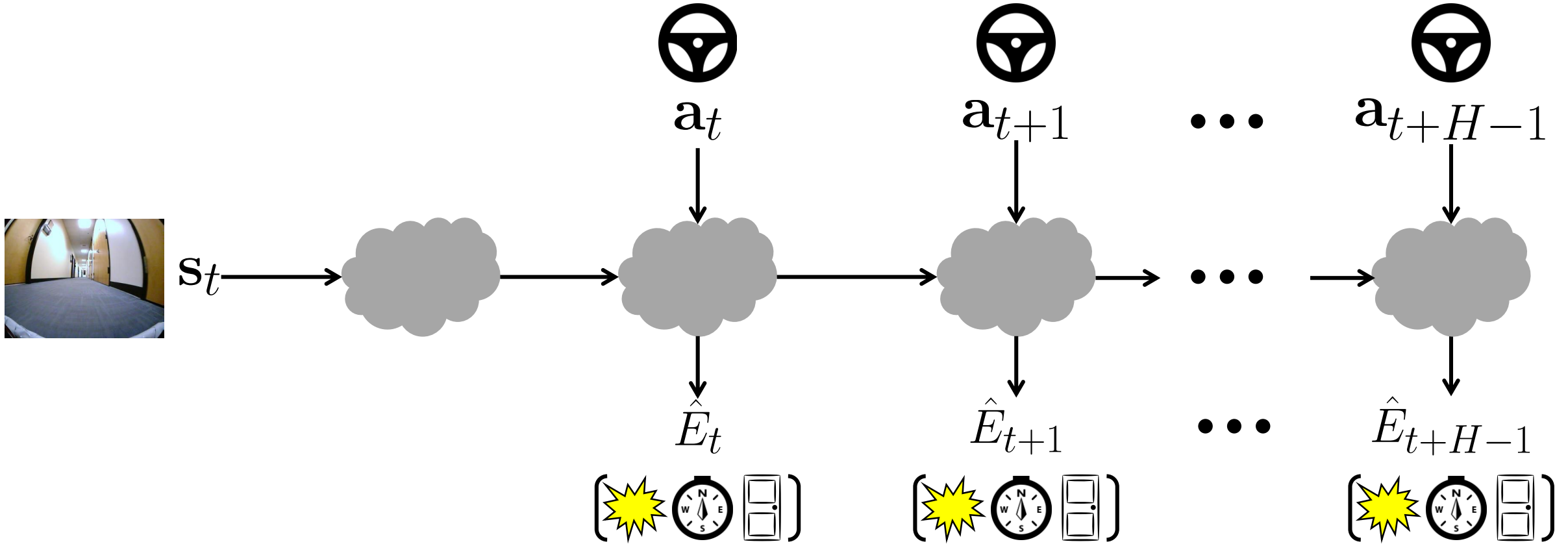
Door Fraction: 0.27



Detect

Predict

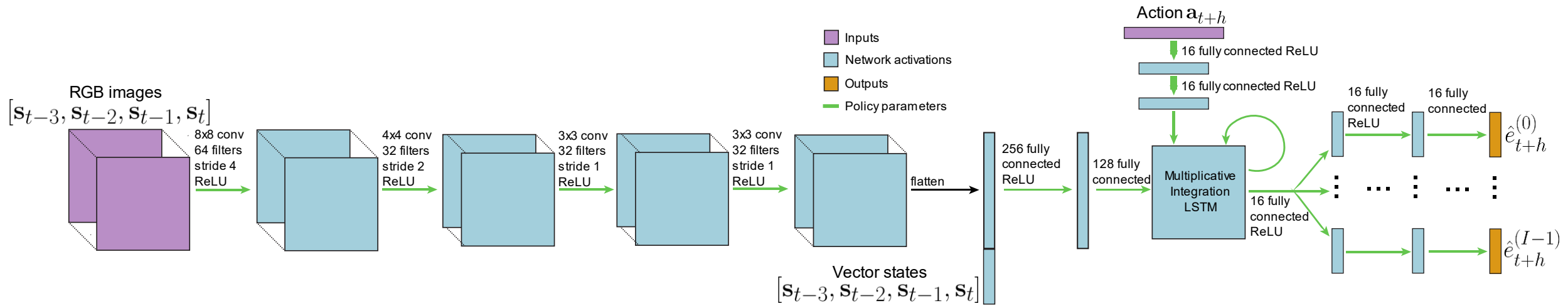
Control



Detect

Predict

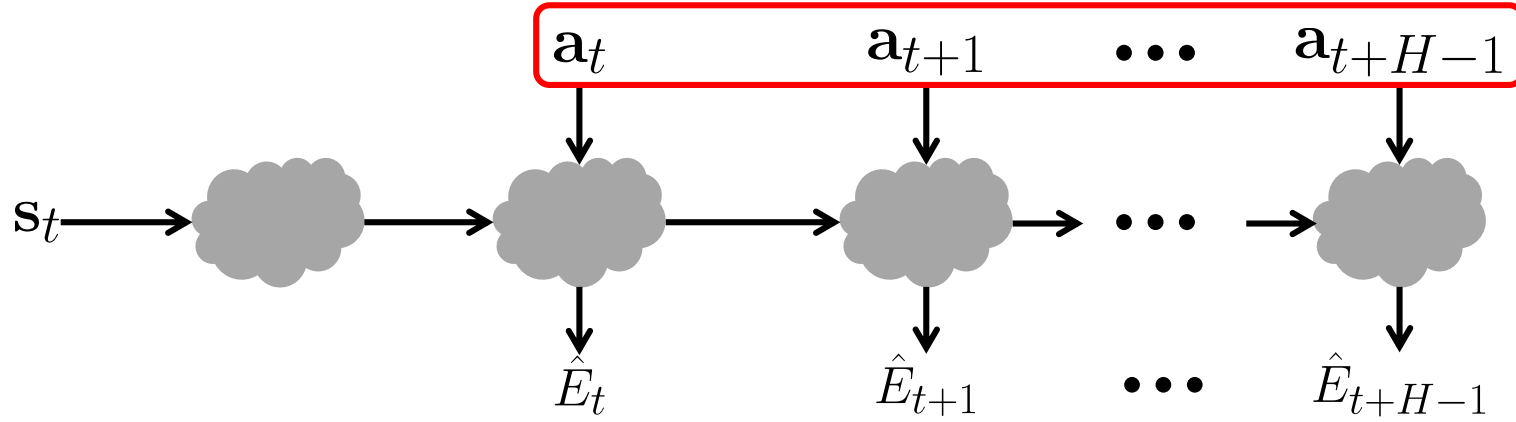
Control



Detect

Predict

Control



$-1 \cdot e_{\text{COLLISION}}$

$1 \cdot e_{\text{HEADING}}$

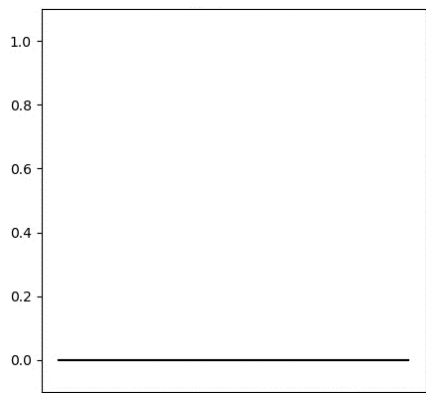
$1 \cdot e_{\text{DOOR}}$

$-100 \cdot e_{\text{COLLISION}} + 1 \cdot e_{\text{HEADING}} + 1 \cdot e_{\text{DOOR}}$

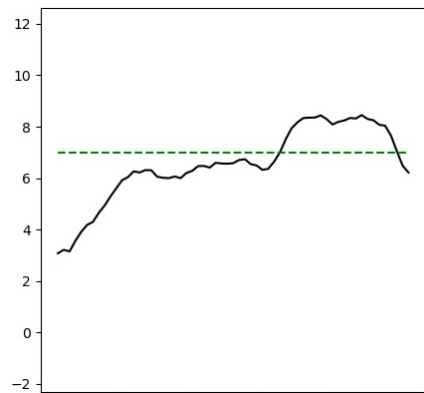




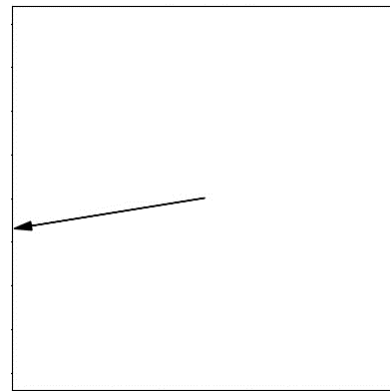
$e_{\text{COLLISION}}$



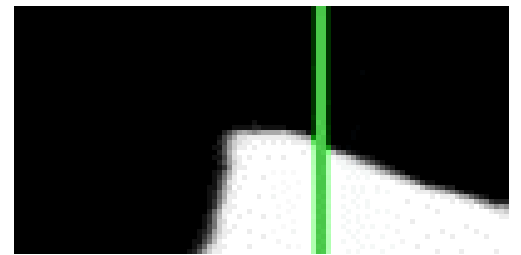
e_{SPEED}



e_{HEADING}



$e_{\text{LANE_SEEN}}$
 $e_{\text{LANE_DIFF}}$



Drive at 7m/s
Avoid collisions

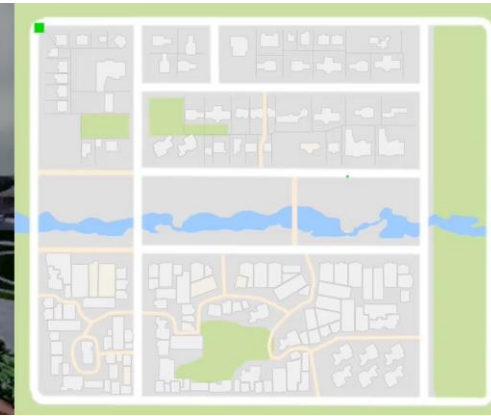
Drive in either lane

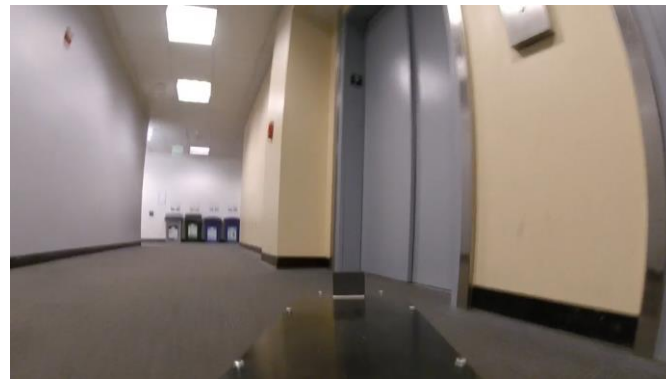


Drive in right lane



CAPs



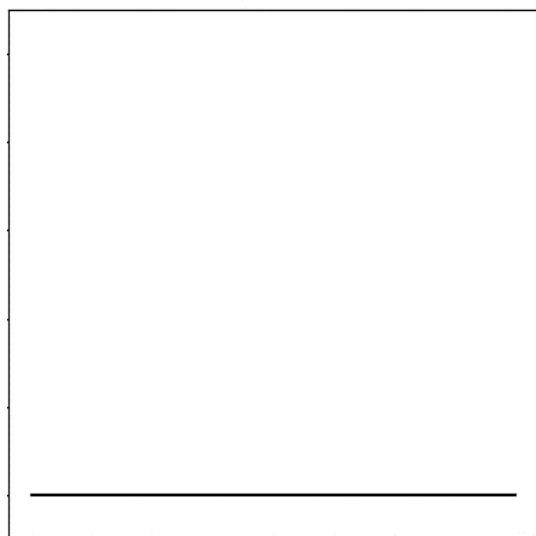


Safety Flexibility

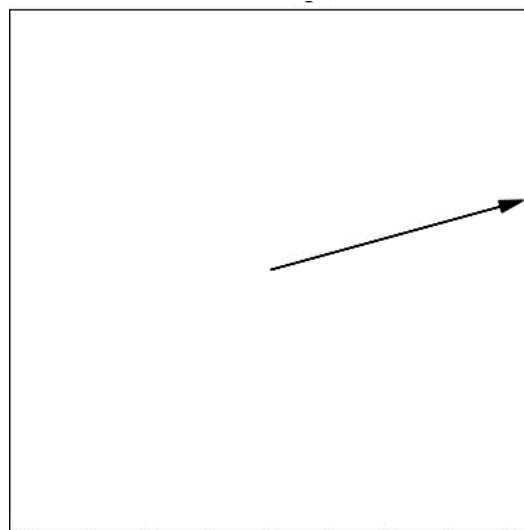
Imitation learning Model-free Model-based



$e_{\text{COLLISION}}$



e_{HEADING}



e_{DOOR}

Door Fraction: 0.27



Collision Avoidance

CAPs



DQL



Avoid collisions
Follow goal heading
Move towards doors

Heading



Flexibility takeaways

- Carefully construct how your policy / model deals with goals
- Model-free methods require extra care to reuse
- Model-based methods are flexible by construction