

Imitation Learning & Behavioral Cloning

Lucas Janson

**CS/Stat 184(0): Introduction to Reinforcement Learning
Fall 2024**

Today

- Feedback from last lecture
- Recap
- Imitation Learning problem statement
- Behavioral Cloning
- DAgger

Feedback from feedback forms

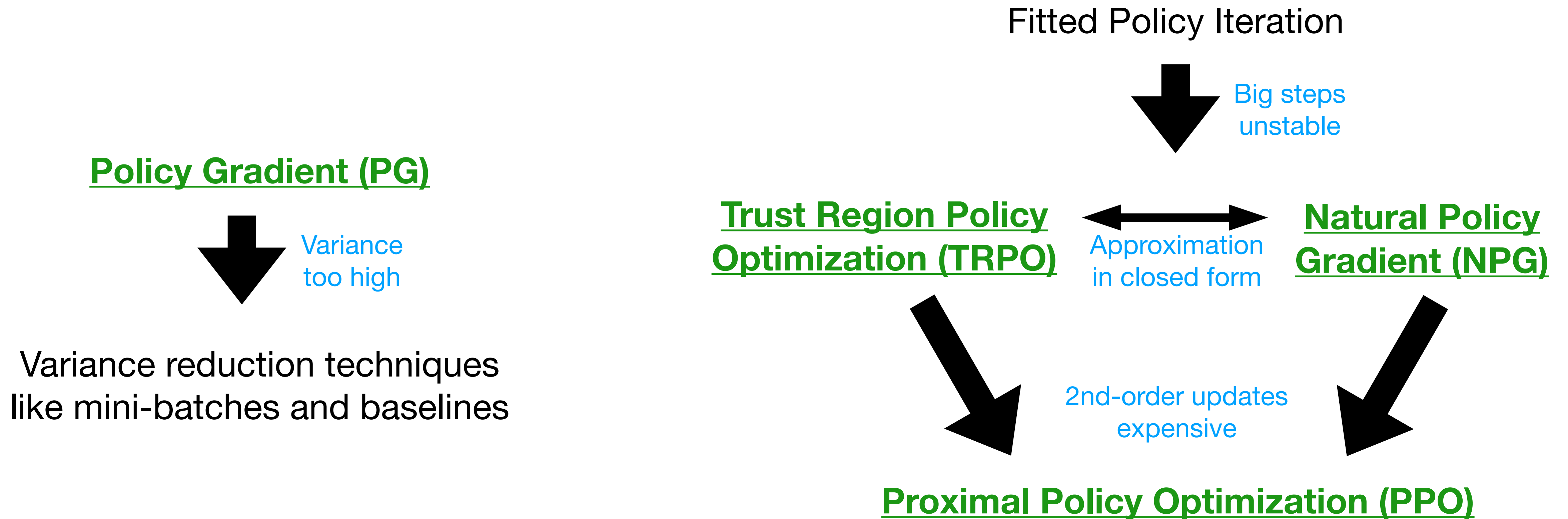
1. Thank you to everyone who filled out the forms!

Today

- ✓ • Feedback from last lecture
- Recap
- Imitation Learning problem statement
- Behavioral Cloning
- DAgger

All Policy Gradient Algorithms in One Slide

Parameterize policy and optimize directly while sampling from MDP



PPO gets 2nd-order optimization benefits over PG and 1st-order computation benefits over TRPO/NPG

“Lack of Exploration” leads to Optimization and Statistical Challenges



Thrun '92

- Suppose $H \approx \text{poly}(|S|)$ & $\mu(s_0) = 1$ (i.e. we start at s_0).
- A randomly initialized policy π^0 has prob. $O(1/3^{|S|})$ of hitting the goal state in a trajectory.
- Thus a sample-based approach, with $\mu(s_0) = 1$, require $O(3^{|S|})$ trajectories.
 - Holds for (sample based) Fitted DP
 - Holds for (sample based) PG/TRPO/NPG/PPO
- Basically, for these approaches, there is no hope of learning the optimal policy if $\mu(s_0) = 1$.

Why not do one trajectory that always moves right?

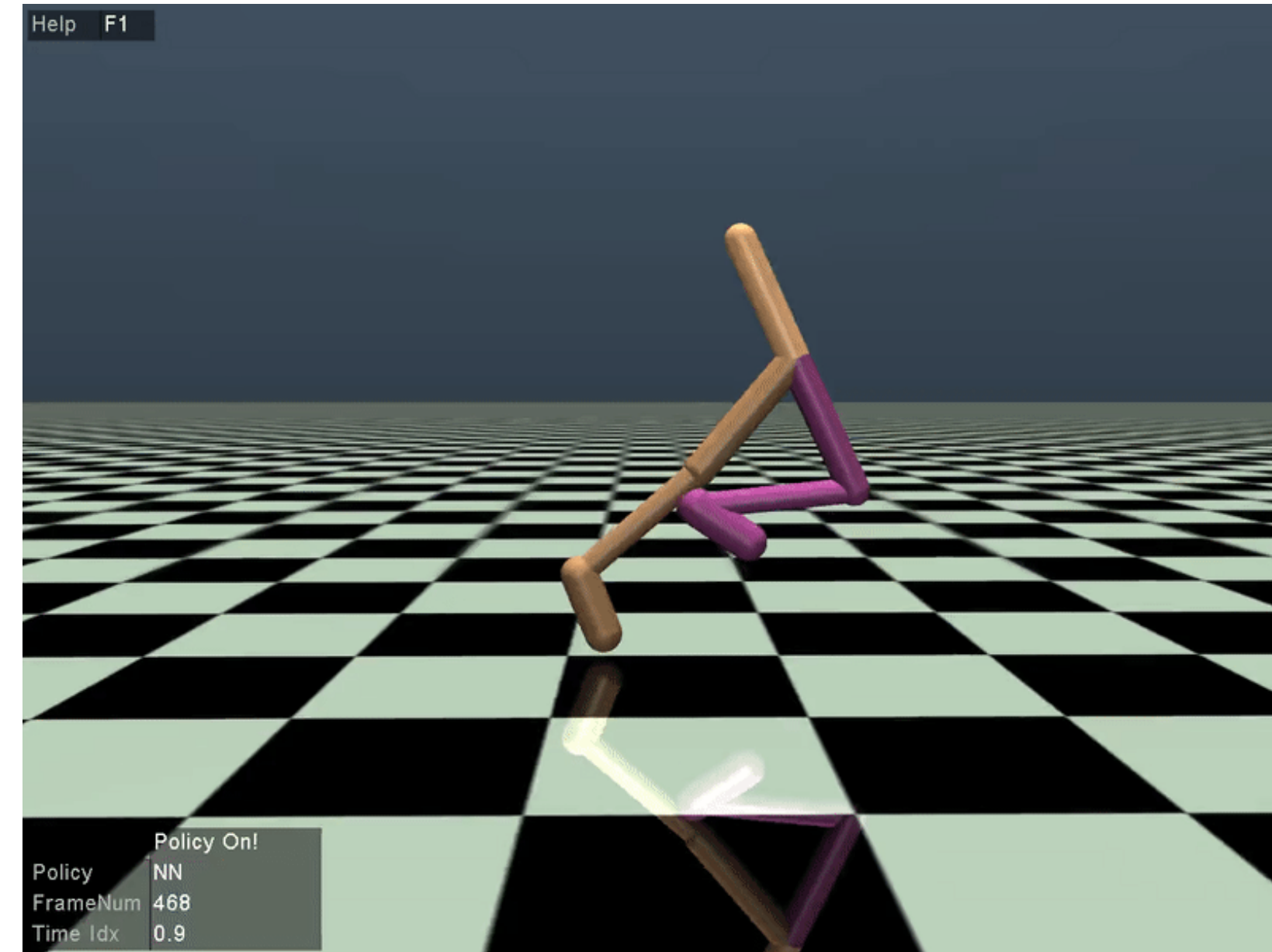
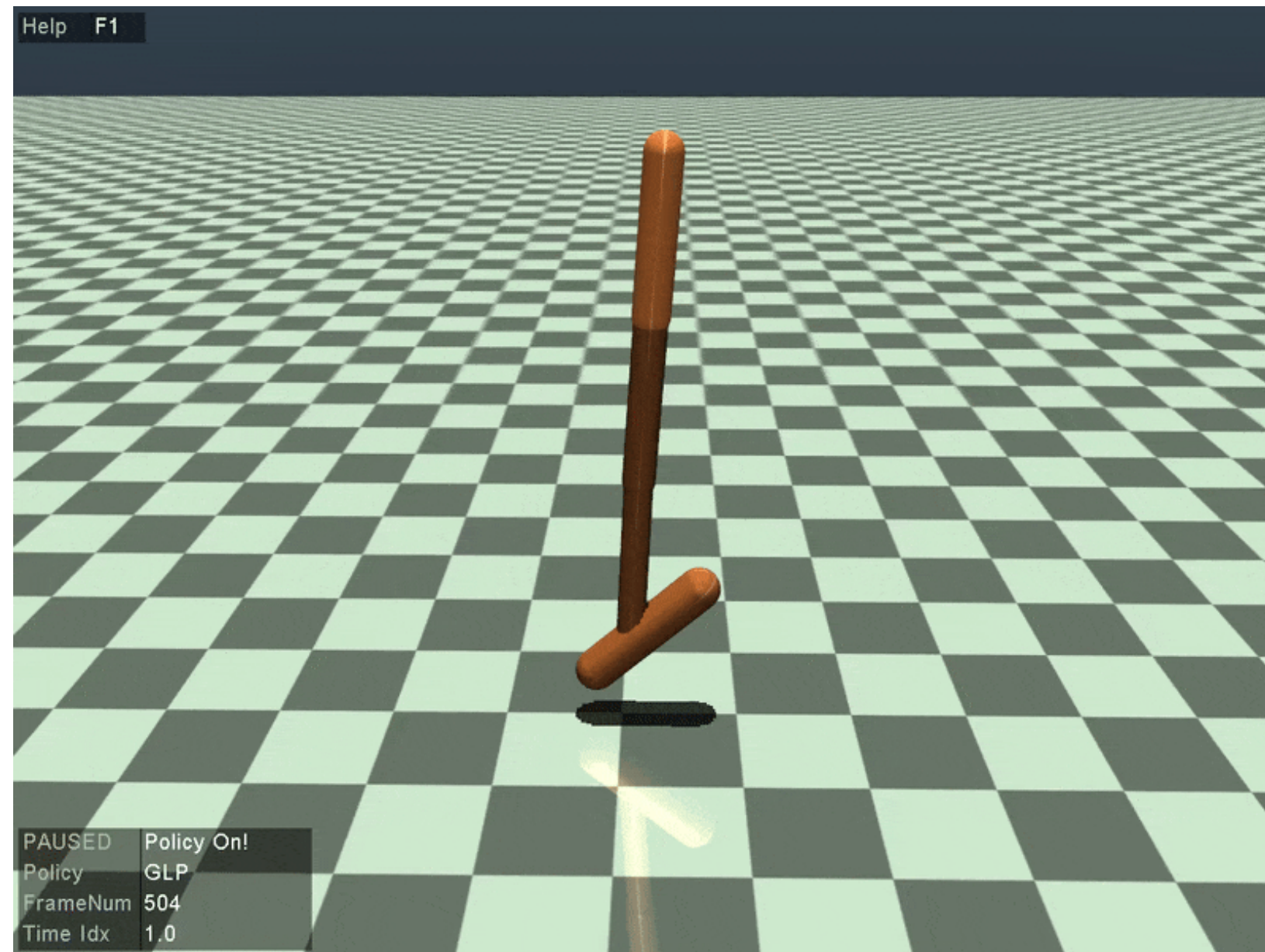
Let's examine the role of μ



Thrun '92

- Suppose that somehow the distribution μ had better coverage.
 - e.g, if μ was uniform overall states in our toy problem, then all approaches we covered would work (with mild assumptions)
 - Theory: **TRPO/NPG/PPO have better guarantees than fitted DP methods** (assuming some “coverage”)
- **Strategies without coverage:**
 - If we have a simulator, sometimes we can **design μ to have better coverage.**
 - this is helpful for robustness as well.
 - **Imitation learning** (next time).
 - An expert gives us samples from a “good” μ .
 - **Explicit exploration:**
 - **UCB-VI:** we'll merge two good ideas!
 - Encourage exploration in PG methods.
 - Try with **reward shaping**

Aside: Brittle policies if we train starting from only from one configuration!

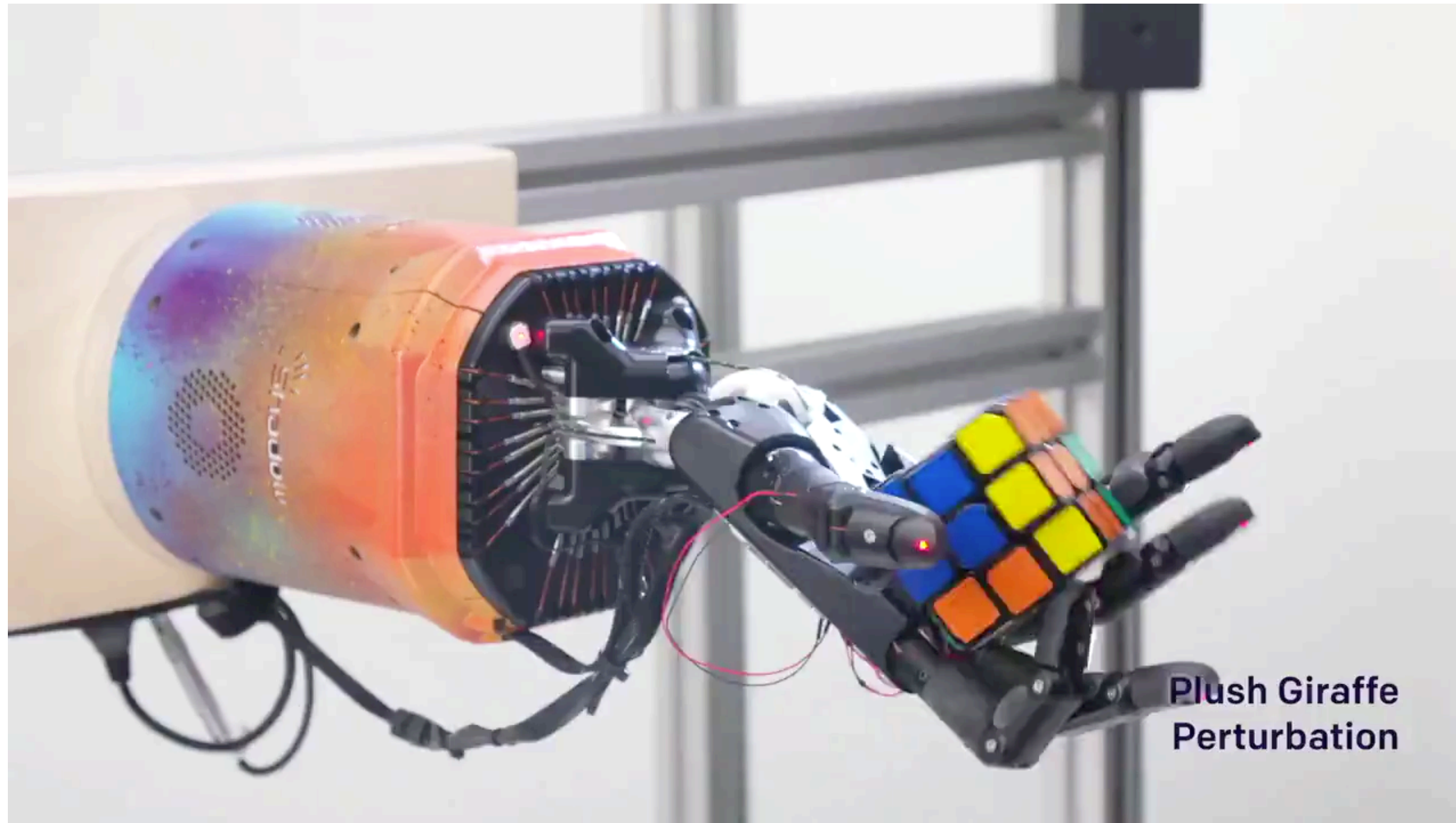


- [Rajeswaran, Lowrey, Todorov, K. 2017]: showed policies optimized for a single starting configuration s_0 are not robust!
- How to fix this?
 - Training from different starting configurations sampled from $s_0 \sim \mu$ fixes this:

$$\max_{\theta} \mathbb{E}_{s_0 \sim \mu} [V^{\theta}(s_0)]$$

Even if starting position concentrated at just one point—good for robustness!

OpenAI: progress on dexterous hand manipulation



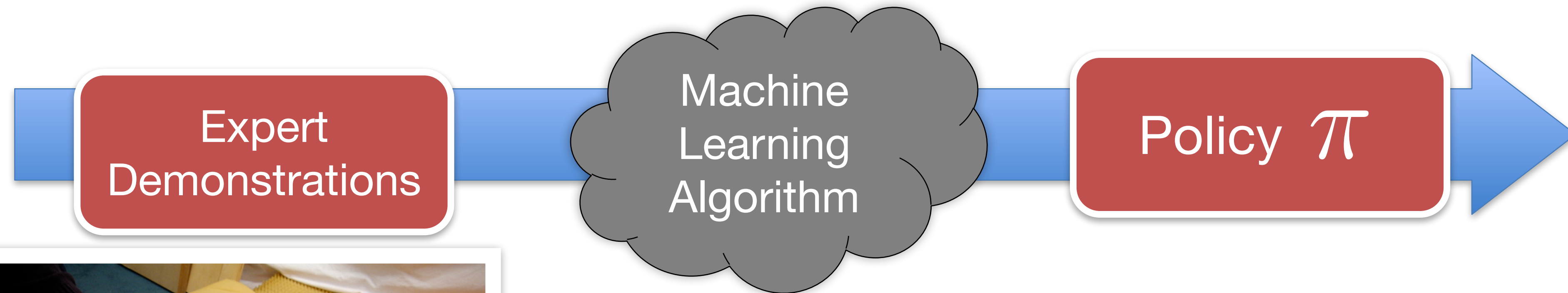
Trained with “domain randomization”

Basically, the measure $s_0 \sim \mu$ was diverse.

Today

- ✓ • Feedback from last lecture
- ✓ • Recap
 - Imitation Learning problem statement
 - Behavioral Cloning
 - DAgger

Imitation Learning



- SVM
- Gaussian Process
- Kernel Estimator
- Deep Networks
- Random Forests
- LWR
- ...

Maps *states* to actions

Learning to Drive by Imitation

[Pomerleau89, Saxena05, Ross11a]

Input:



Camera Image



Output:

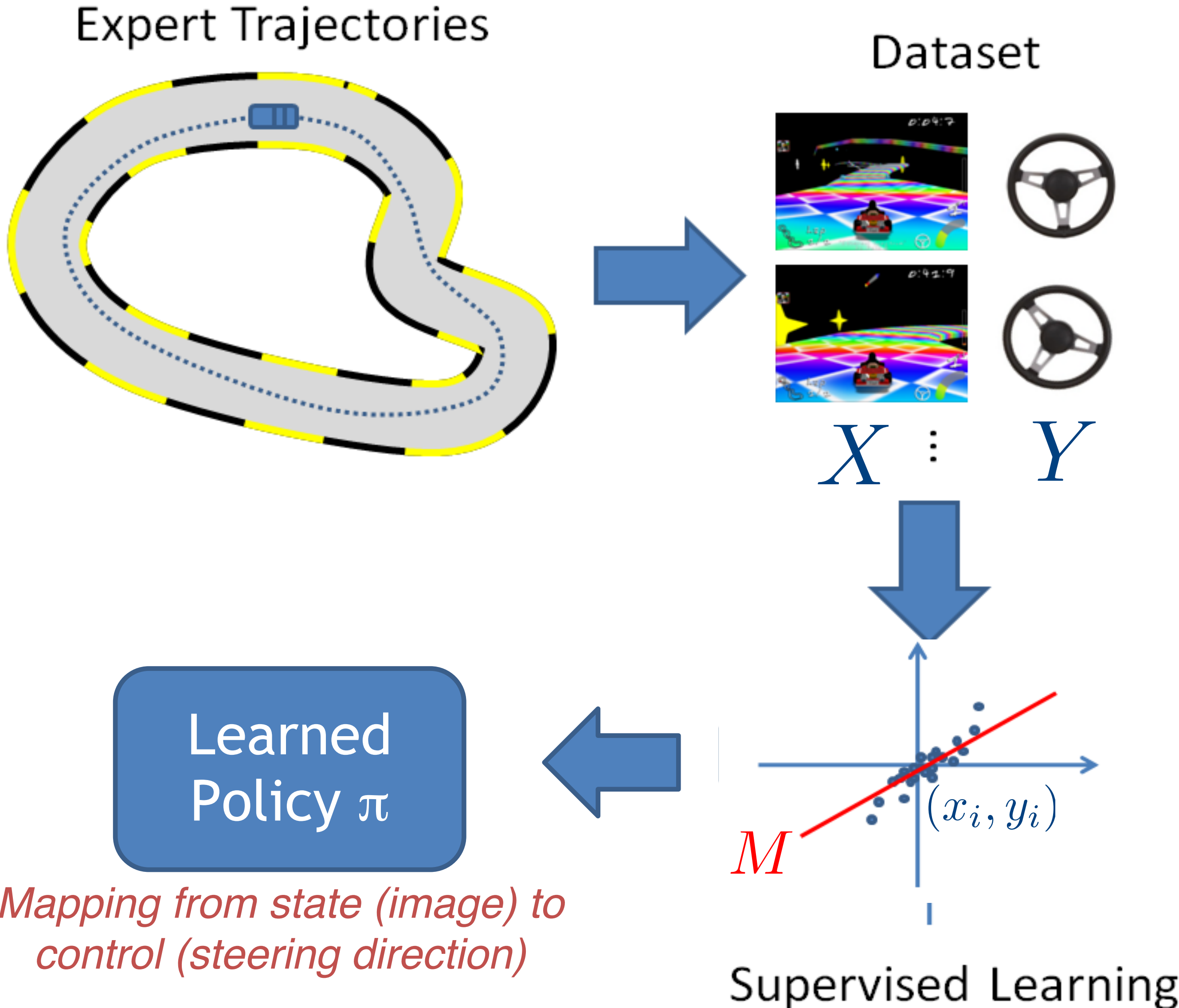


Steering Angle
in $[-1, 1]$

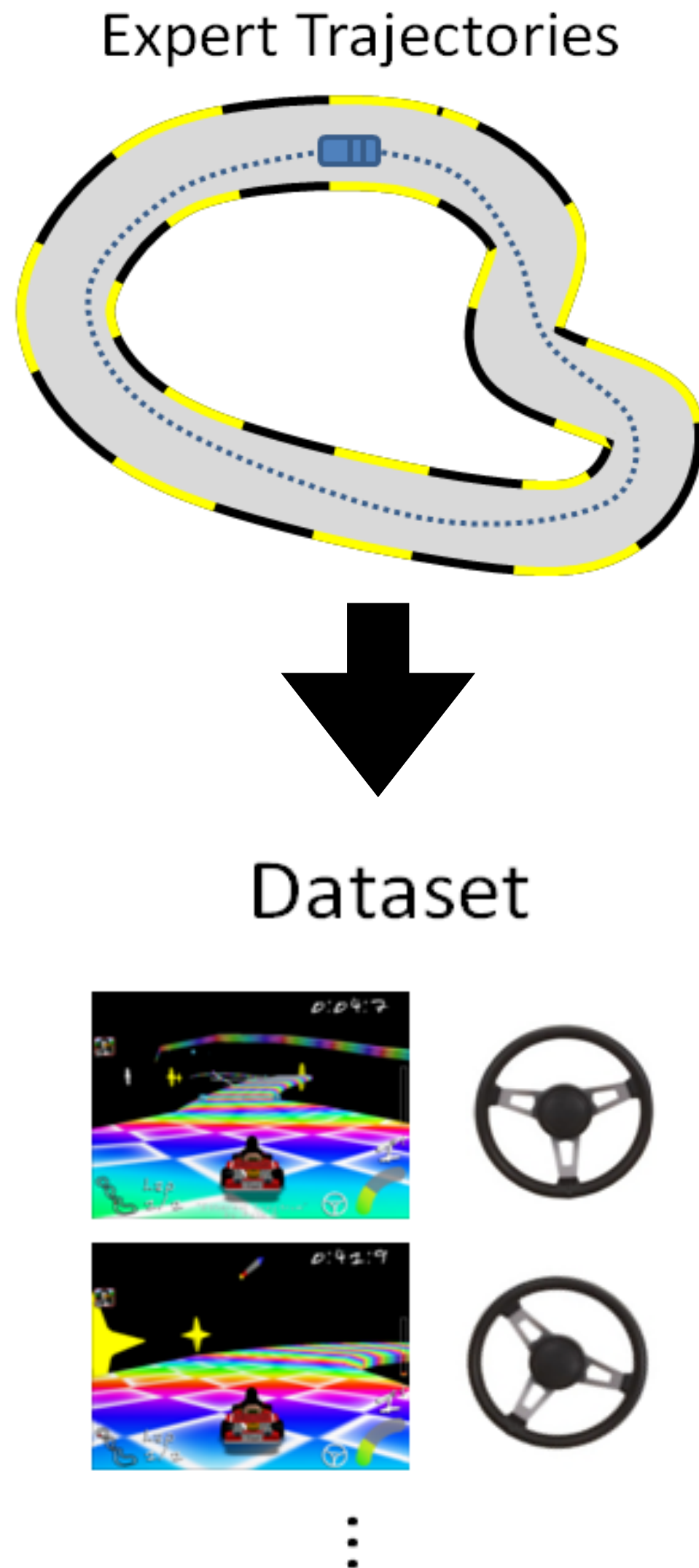
Today

- ✓ • Feedback from last lecture
- ✓ • Recap
- ✓ • Imitation Learning problem statement
 - Behavioral Cloning
 - DAgger

Supervised Learning Approach: Behavior Cloning



Let's formalize the offline IL Setting and the Behavior Cloning algorithm



Finite horizon MDP \mathcal{M}

Ground truth reward $r(s, a) \in [0, 1]$ is unknown;
Assume the expert has a good policy π^* (not necessarily opt)

We have a dataset of M trajectories: $\mathcal{D} = \{\tau_1, \dots, \tau_M\}$,
where $\tau_i = (s_h^i, a_h^i)_{h=0}^{H-1} \sim \rho_{\pi^*}$

Goal: learn a policy from \mathcal{D} that is as good as the expert π^*

Let's formalize the Behavior Cloning (BC) algorithm

BC Algorithm input: a restricted policy class $\Pi = \{\pi : S \mapsto \Delta(A)\}$

BC is a Reduction to Supervised Learning:

$$\hat{\pi} = \arg \min_{\pi \in \Pi} \sum_{i=1}^M \sum_{h=0}^{H-1} \ell(\pi, s_h^i, a_h^i)$$

$\ell(\pi, s, a)$ is a loss function with many choices:

1. Classification (0/1) loss: $\mathbf{1}[\pi(s) \neq a]$
2. Negative log-likelihood (NLL): $\ell(\pi, s, a) = -\ln \pi(a | s)$
3. square loss (i.e., regression for continuous action): $\ell(\pi, s, a) = \|\pi(s) - a\|_2^2$

Theorem: IL is (almost) as easy as SL

$$\hat{\pi} = \arg \min_{\pi \in \Pi} \sum_{i=1}^M \sum_{h=0}^{H-1} \ell(\pi, s_h^i, a_h^i)$$

Note a training and testing “mismatch”

Theorem [BC Performance]:

suppose we assume supervised learning succeeds, with ϵ classification error:

$$\mathbb{E}_{\tau \sim \rho_{\pi^*}} \left[\frac{1}{H} \sum_{h=0}^{H-1} \mathbf{1} [\hat{\pi}(s_h) \neq \pi^*(s_h)] \right] \leq \epsilon,$$

(where π^* is the expert policy, which need not be optimal)

then we have:

$$|V^{\pi^*} - V^{\hat{\pi}}| \leq H^2 \epsilon$$

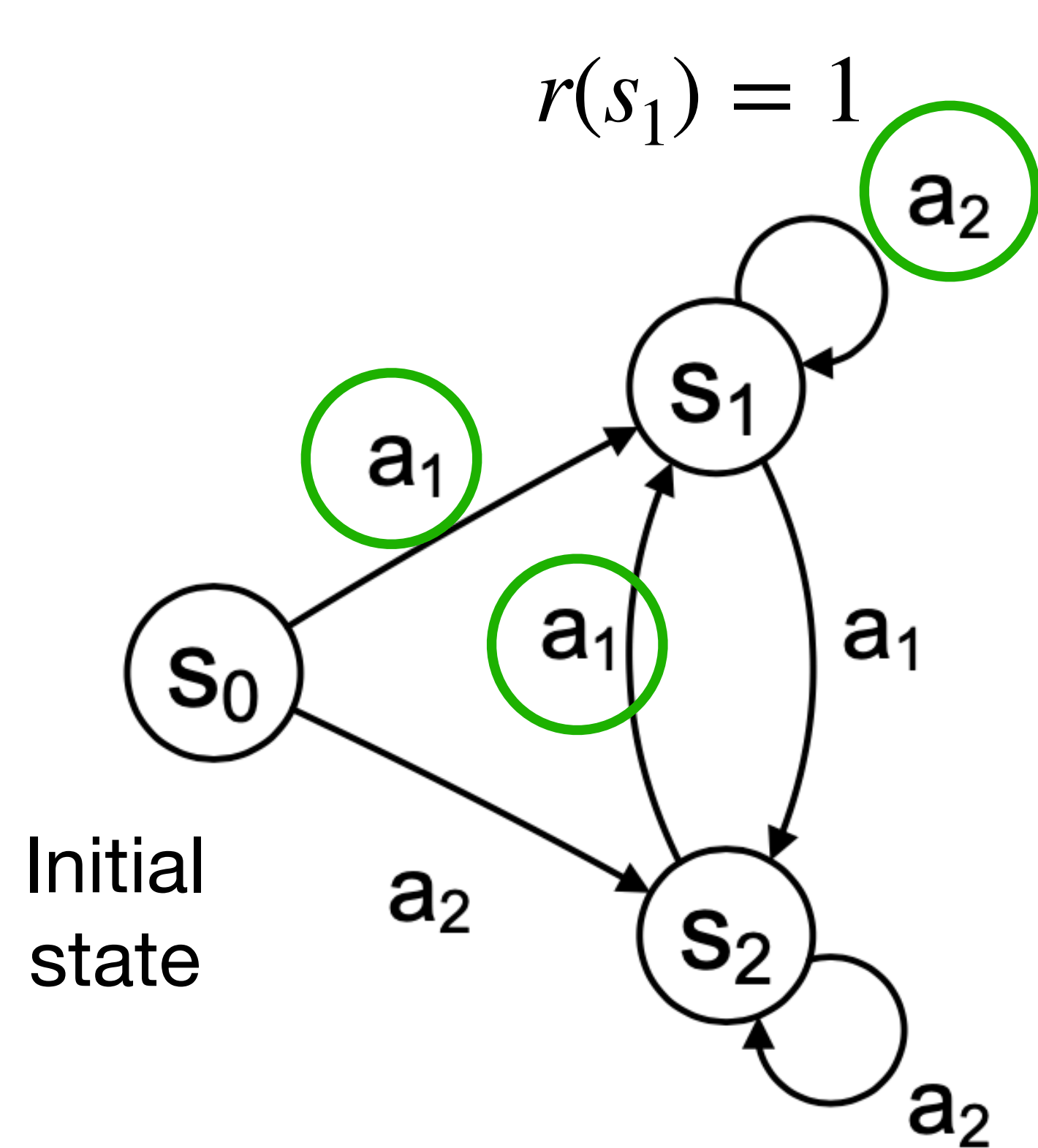
The quadratic amplification is annoying

Proof:

By the PDL

$$\begin{aligned} |V^{\pi^*}(s) - V^{\hat{\pi}}(s)| &= \left| \mathbb{E}_{\tau \sim \rho_{\pi^*}} \left[\sum_{h=0}^{H-1} A_h^{\hat{\pi}}(s_h, a_h) \right] \right| \\ &= \left| \mathbb{E}_{s_1, \dots, s_H \sim \rho_{\pi^*}} \left[\sum_{h=0}^{H-1} A_h^{\hat{\pi}}(s_h, \pi^*(s_h)) \right] \right| \\ &\leq H \left| \mathbb{E}_{\tau \sim \rho_{\pi^*}} \left[\sum_{h=0}^{H-1} \mathbf{1}[\hat{\pi}(s_h) \neq \pi^*(s_h)] \right] \right| \\ &\leq H^2 \epsilon \end{aligned}$$

Distribution Shift Example (H^2 factor is tight)



Assume SL returns the policy $\hat{\pi}$:

$$\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w/ prob } 1 - H\epsilon \\ a_2 & \text{w/ prob } H\epsilon \end{cases}, \quad \hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

This policy has good supervised learning error:

$$\mathbb{E}_{\tau \sim \rho_{\pi^*}} \left[\frac{1}{H} \sum_{h=0}^{H-1} \mathbf{1} [\hat{\pi}(s_h) \neq \pi^*(s_h)] \right] = \epsilon$$

note: while $\hat{\pi}(s_2) \neq \pi^*(s_2)$, state s_2 is never visited under π^*

We have **quadratic degradation** (in H):

$$V_0^{\hat{\pi}}(s_0) = (1 - H\epsilon) \cdot V_0^{\pi^*}(s_0) + H\epsilon \cdot 0 = V_0^{\pi^*}(s_0) - \epsilon H(H - 1)$$

Intuition: once we make a mistake at s_0 , we end up in s_2 which is not in the training data!

Opt policy:

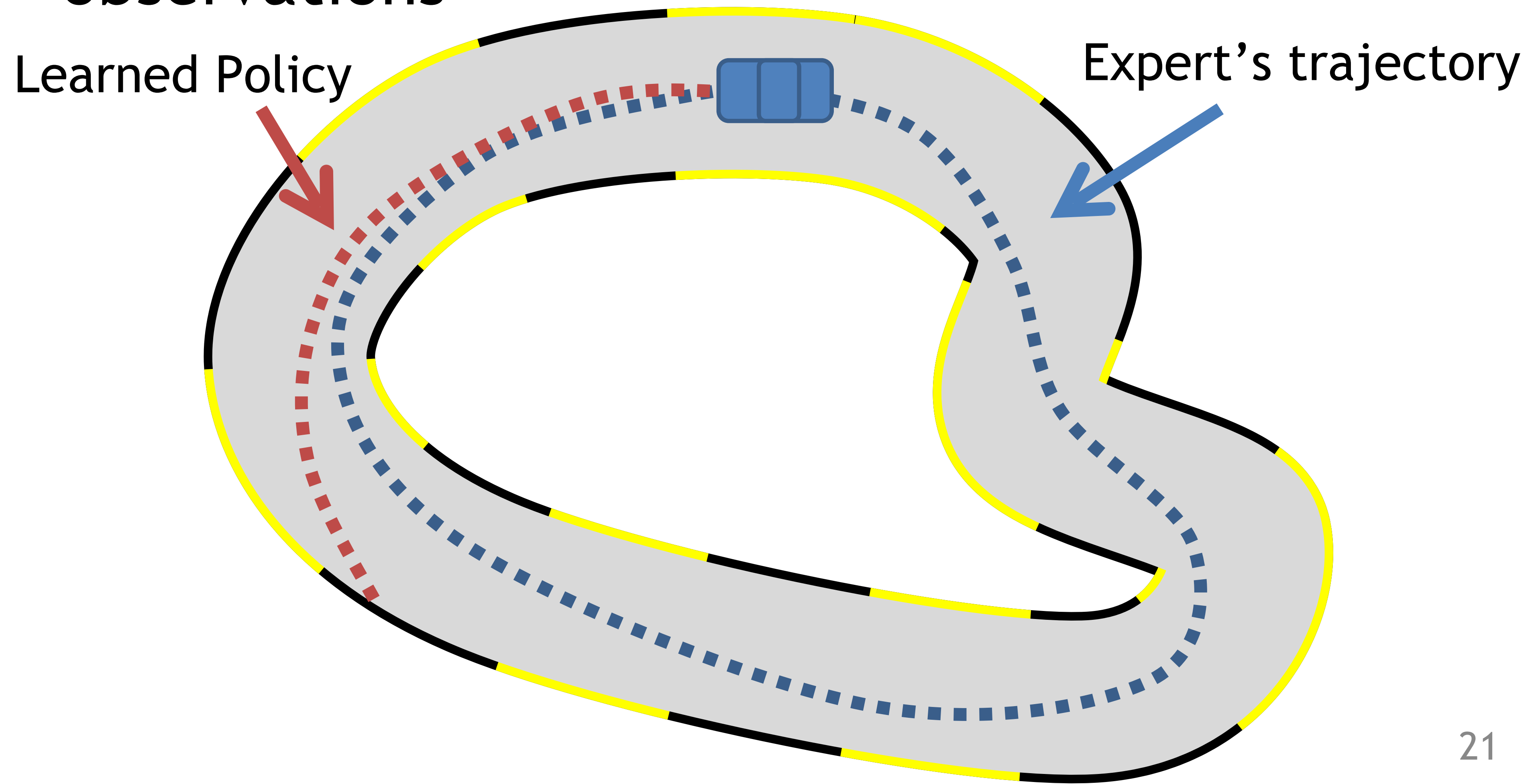
Under ρ_{π^*} , trajectory is s_0, s_1, s_1, \dots

$$\rho_{\pi^*}(s_h = s_2) = 0$$

$$V_0^{\pi^*}(s_0) = H - 1$$

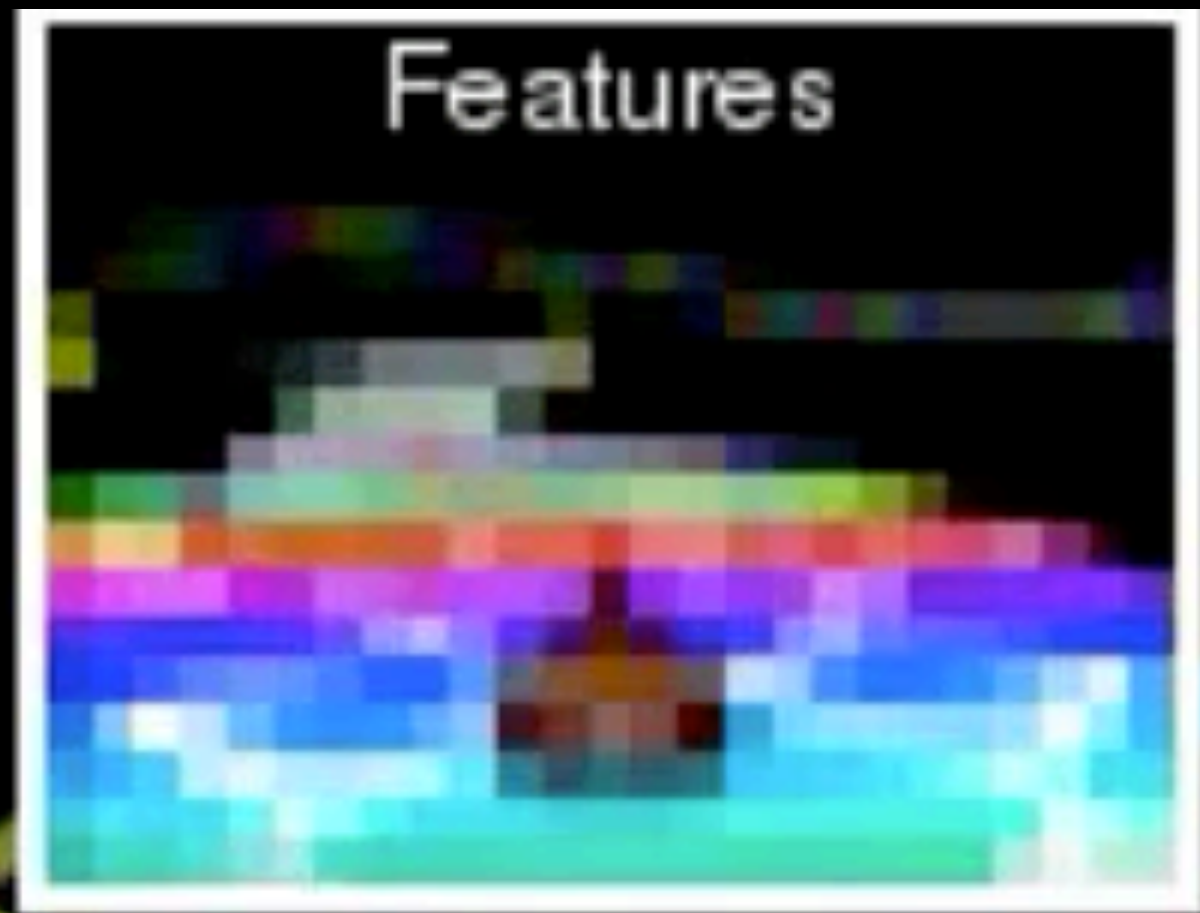
What could go wrong?

- Predictions affect future inputs/ observations



Expert Demos

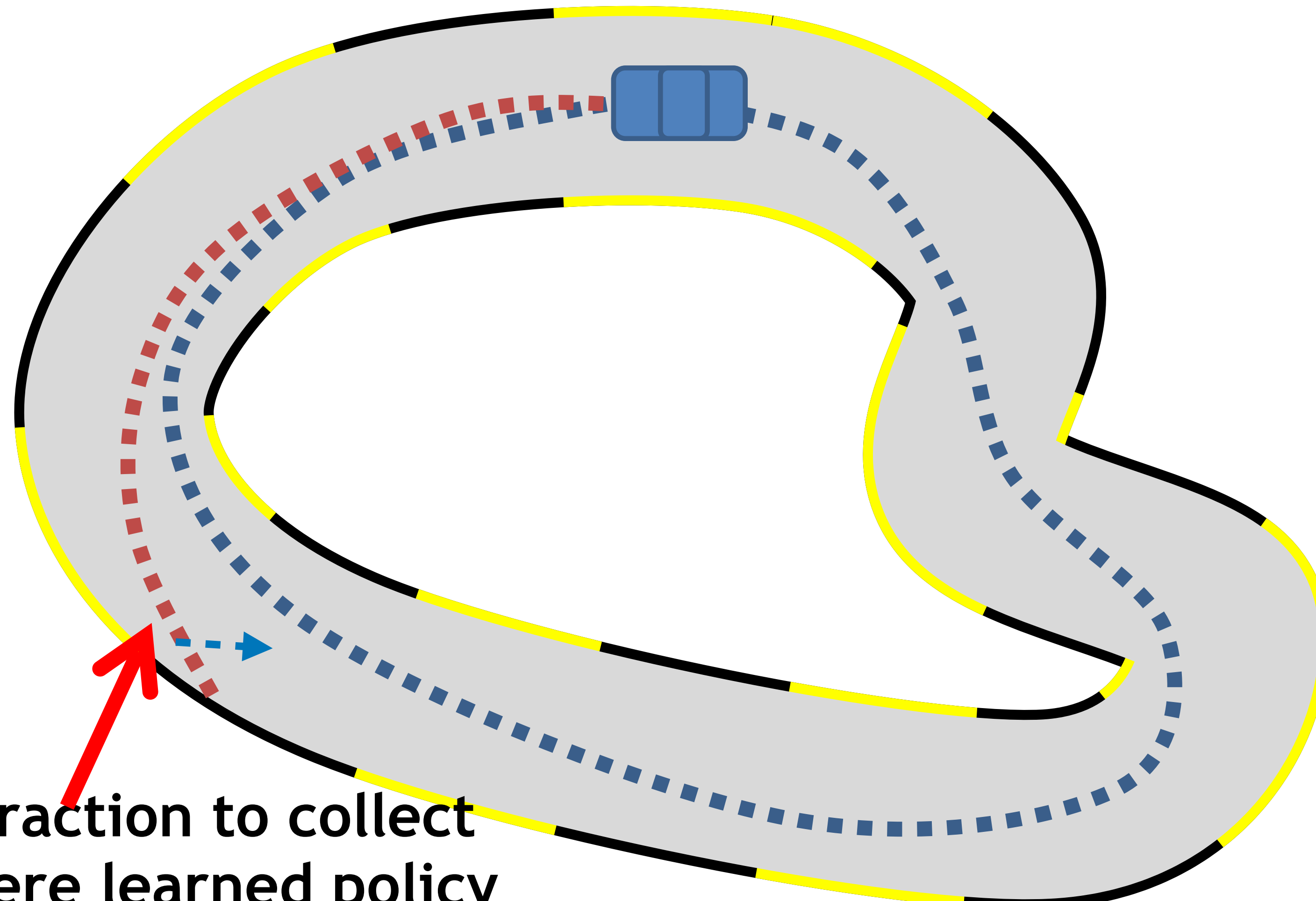




Today

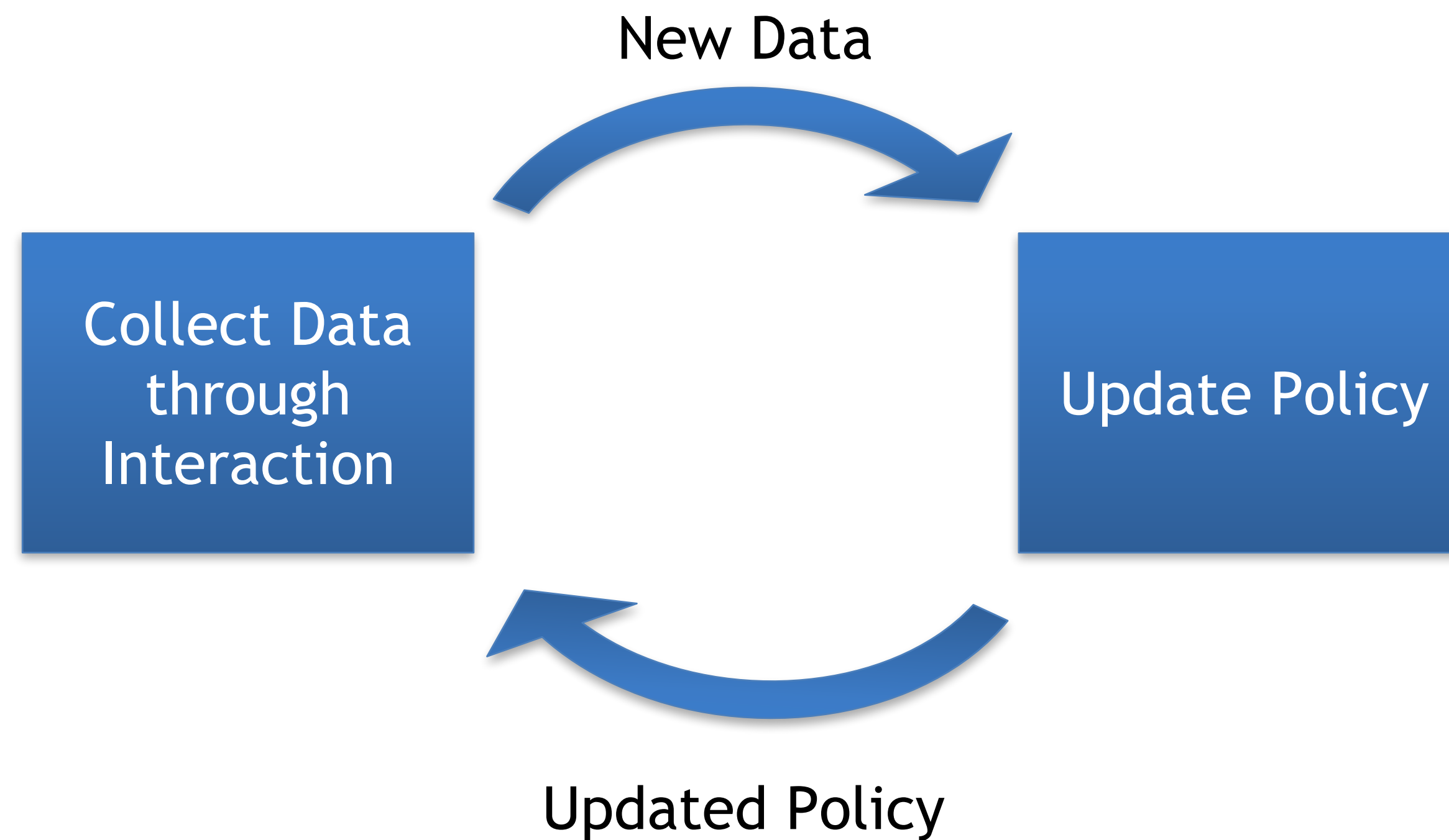
- ✓ • Feedback from last lecture
- ✓ • Recap
- ✓ • Imitation Learning problem statement
- ✓ • Behavioral Cloning
 - DAgger

Intuitive solution: Interaction



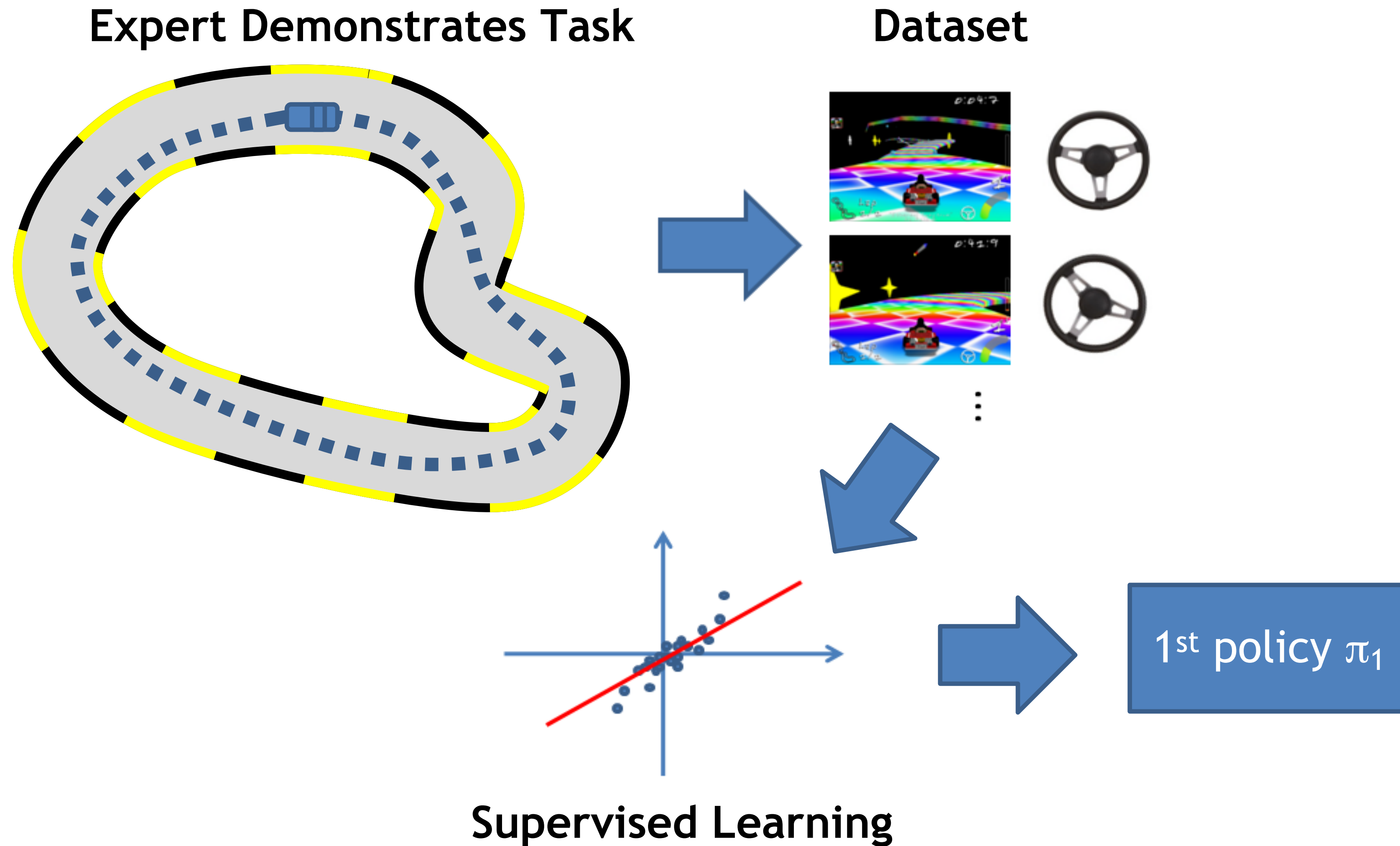
Use interaction to collect data where learned policy goes

General Idea: Iterative Interactive Approach



Dagger: Dataset Aggregation ^[Ross11a]

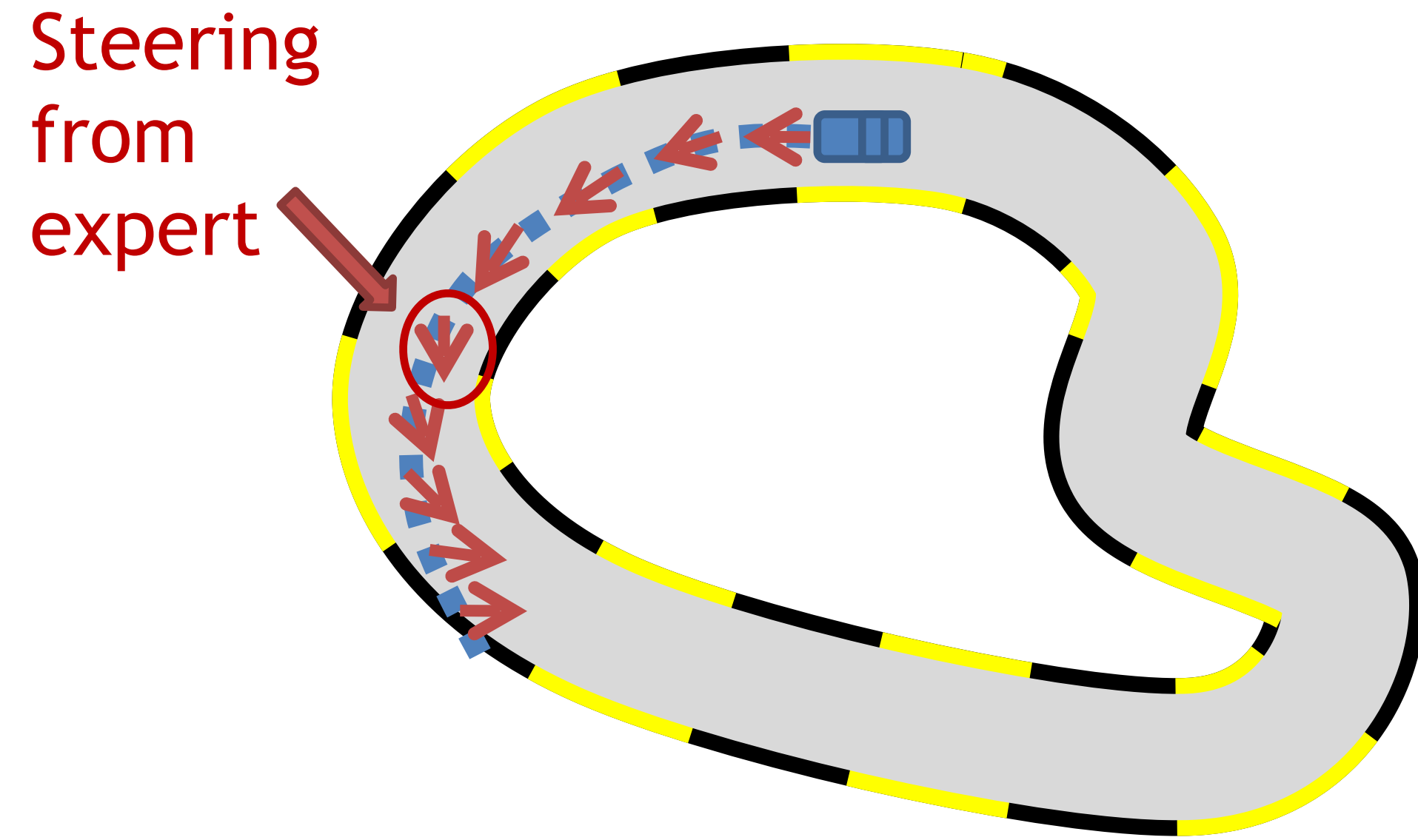
0th iteration



Dagger: Dataset Aggregation ^[Ross11a]

1st iteration

Execute π_1 and Query Expert

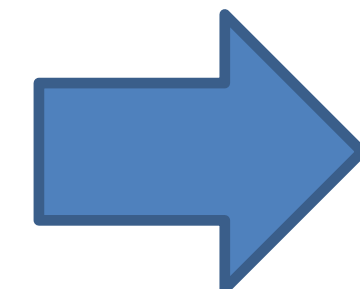
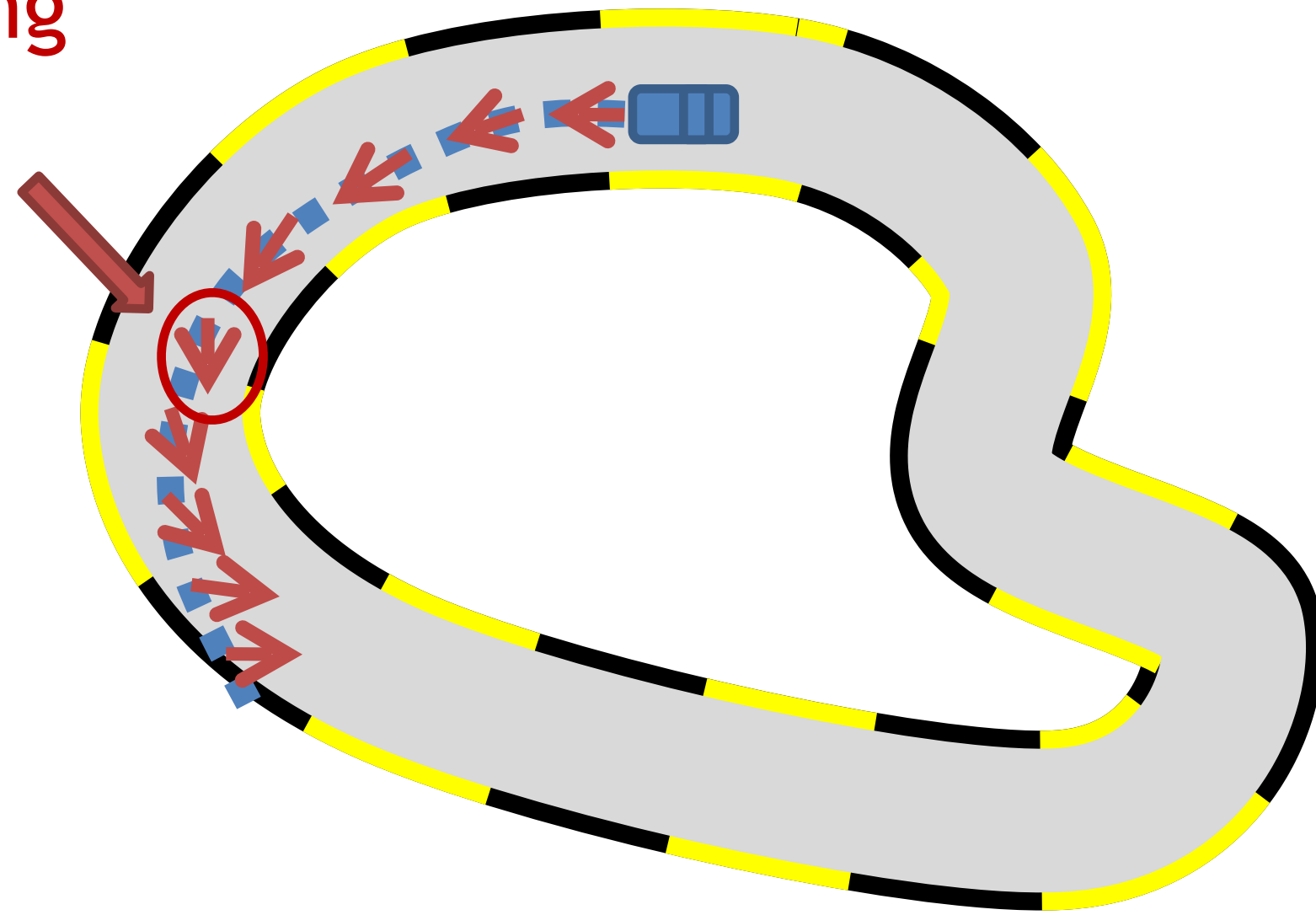


Dagger: Dataset Aggregation ^[Ross11a]

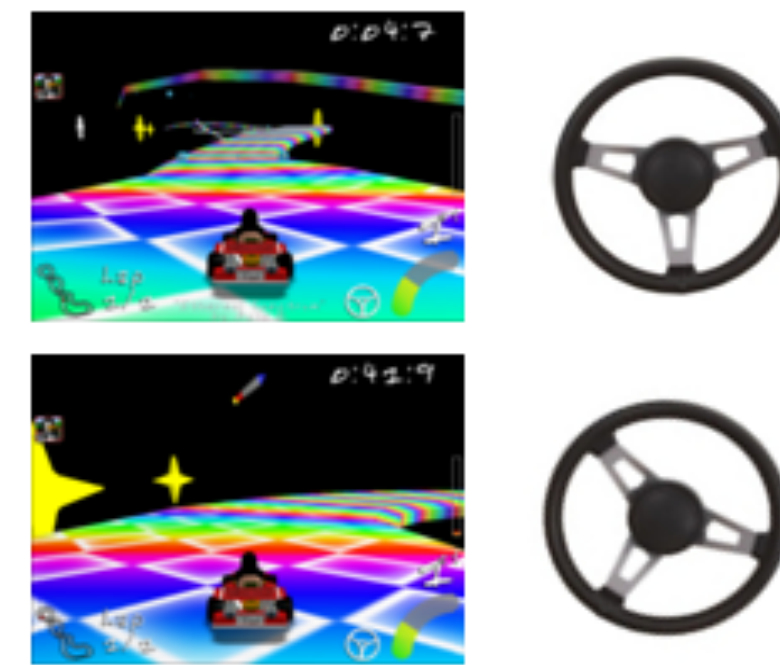
1st iteration

Execute π_1 and Query Expert

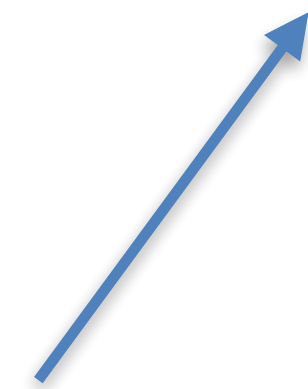
Steering
from
expert



New Data



States from
the learned policy

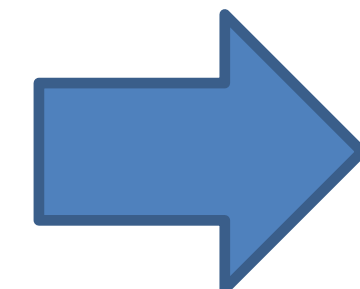
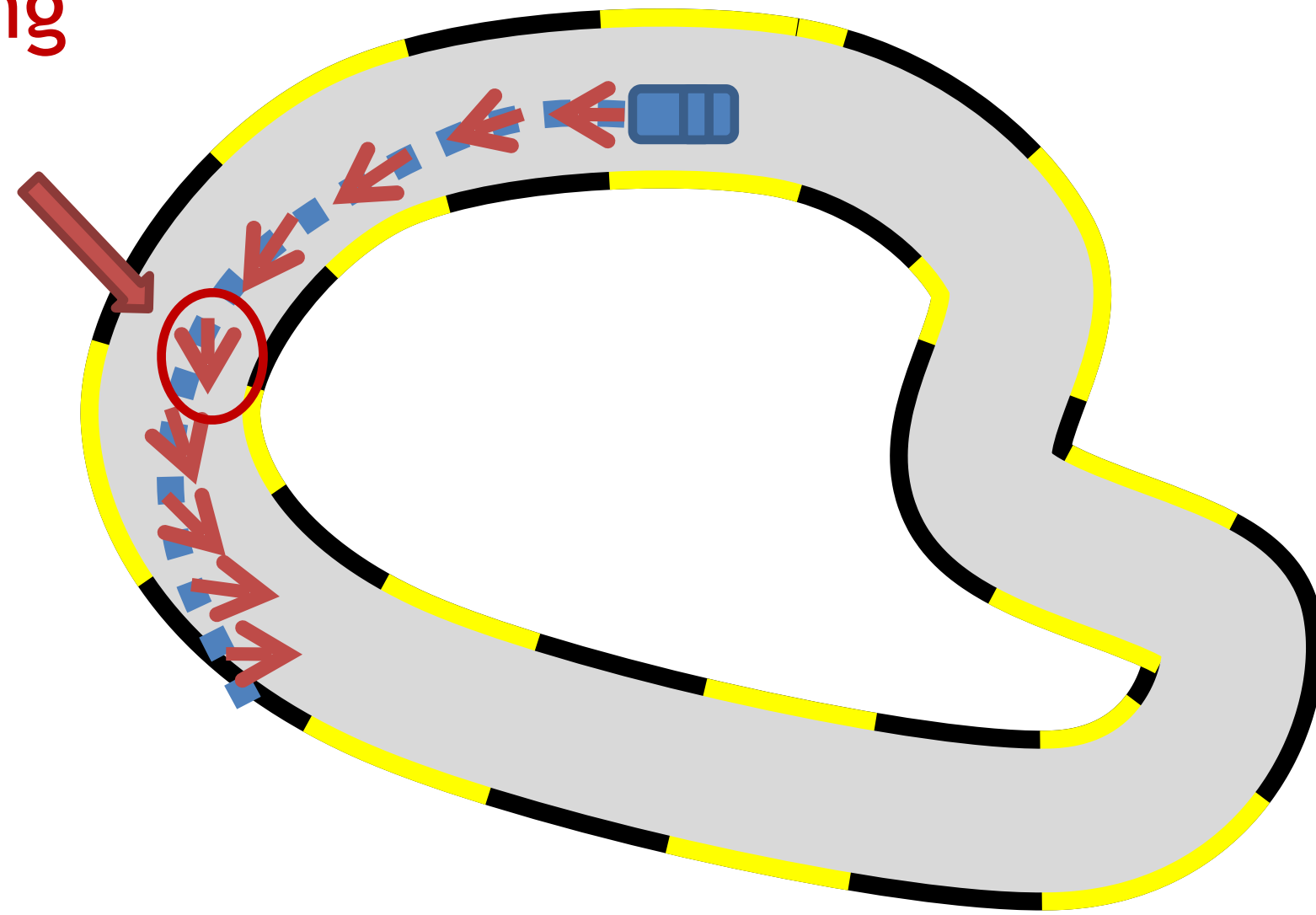


Dagger: Dataset Aggregation ^[Ross11a]

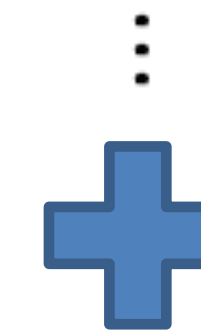
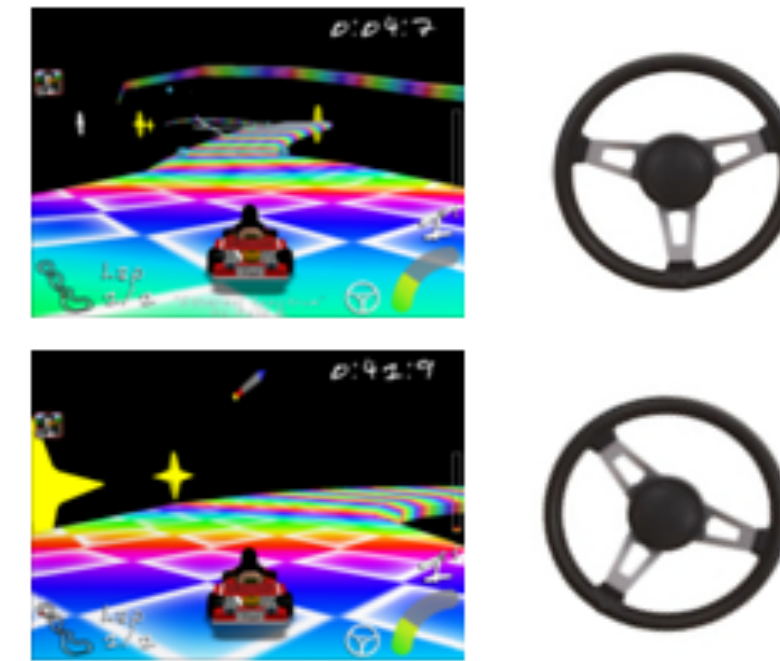
1st iteration

Execute π_1 and Query Expert

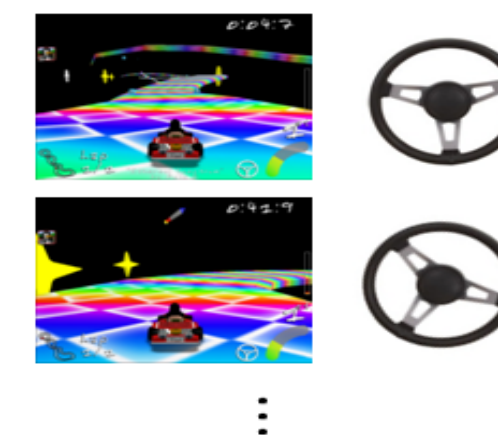
Steering
from
expert



New Data

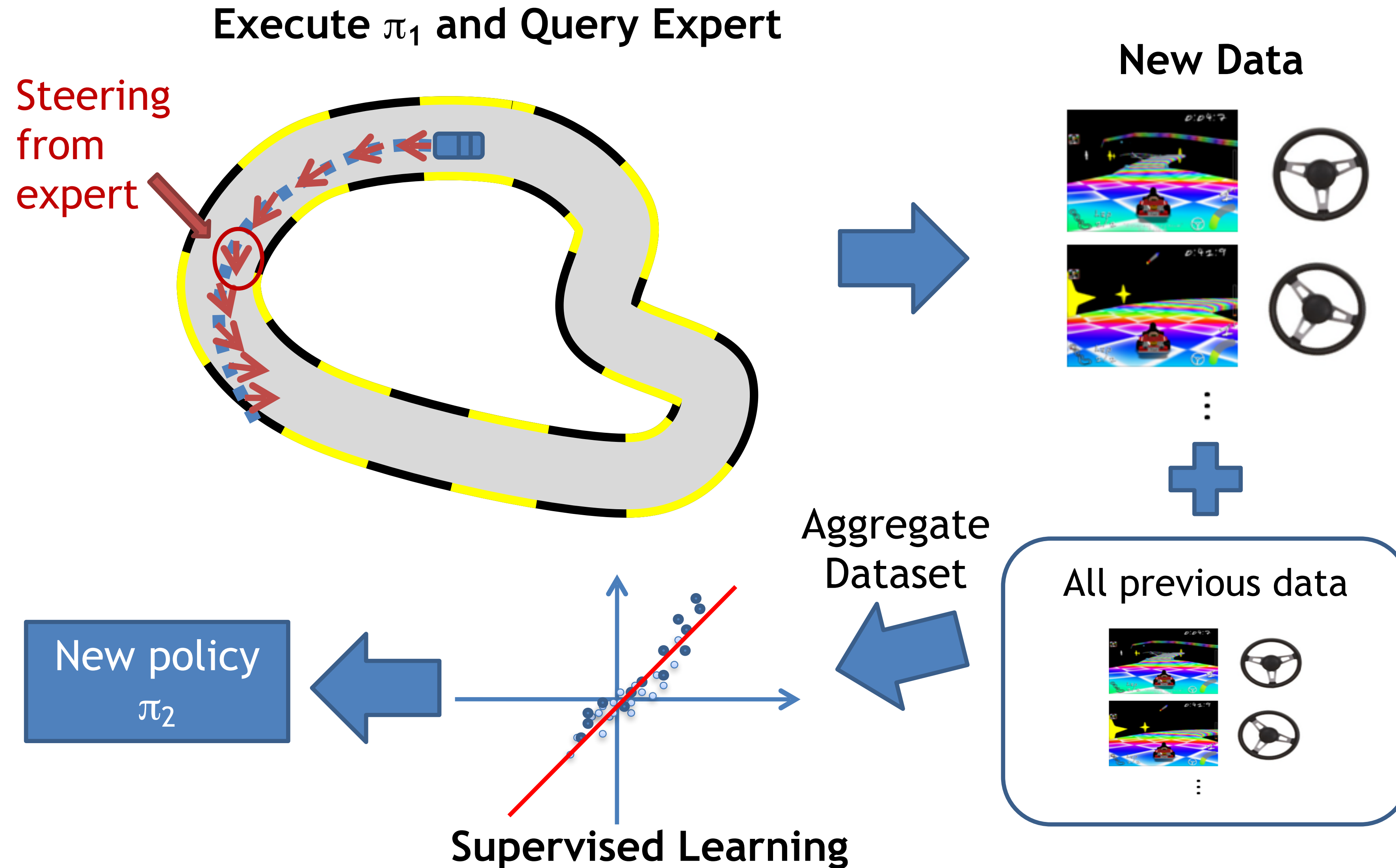


All previous data



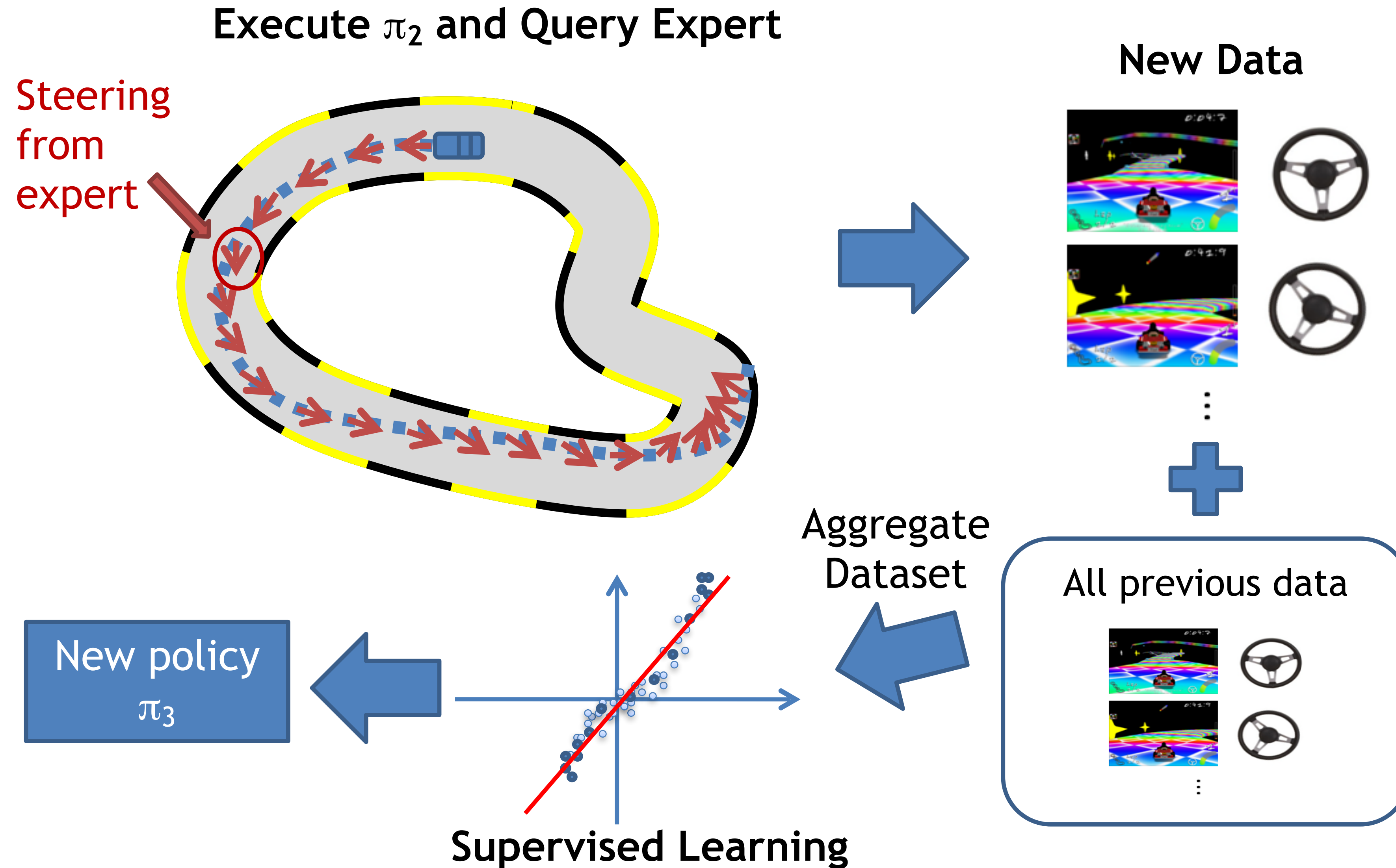
Dagger: Dataset Aggregation [Ross11a]

1st iteration



Dagger: Dataset Aggregation [Ross11a]

2nd iteration

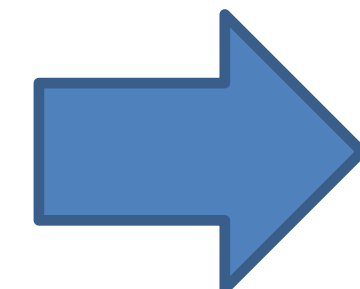
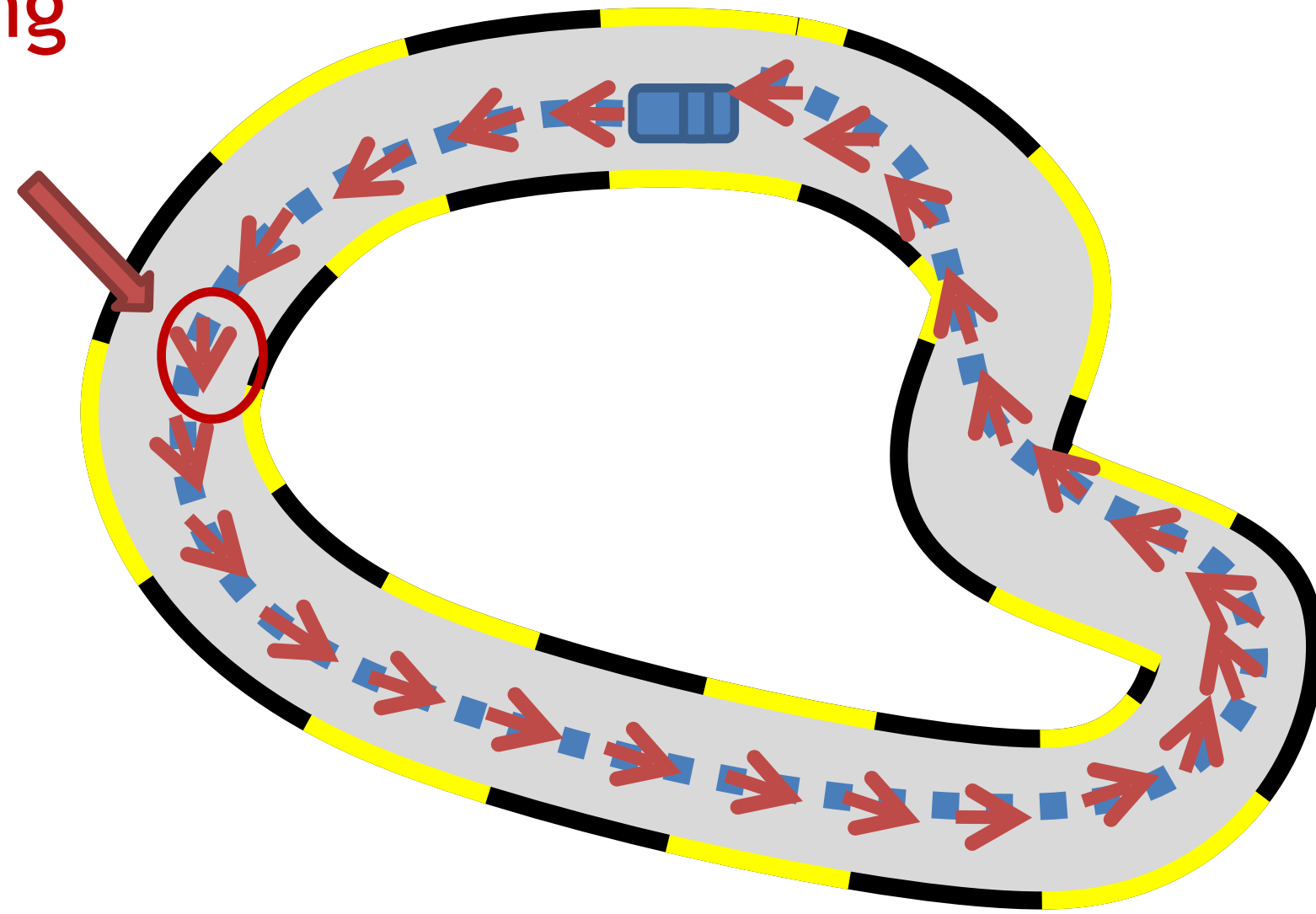


Dagger: Dataset Aggregation [Ross11a]

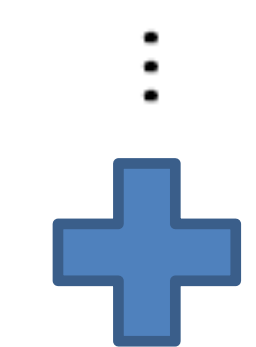
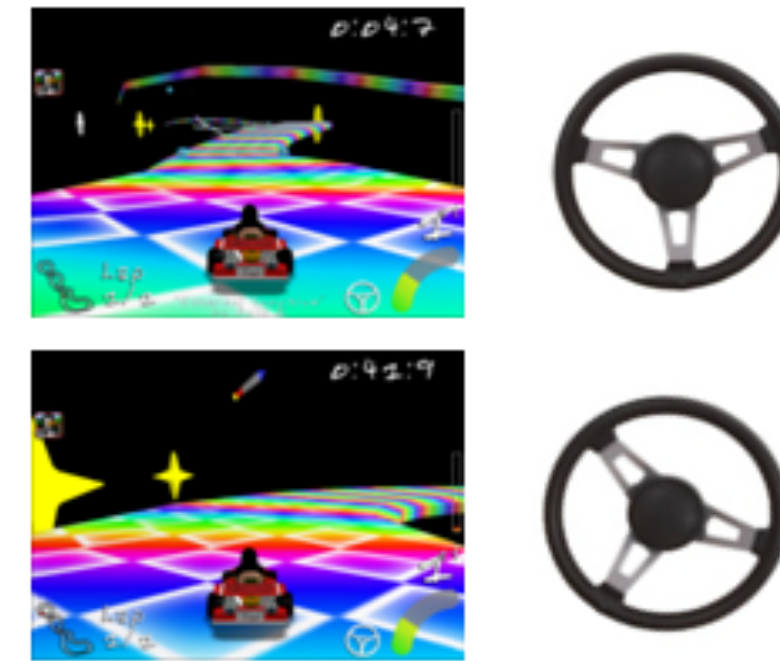
n^{th} iteration

Execute π_{n-1} and Query Expert

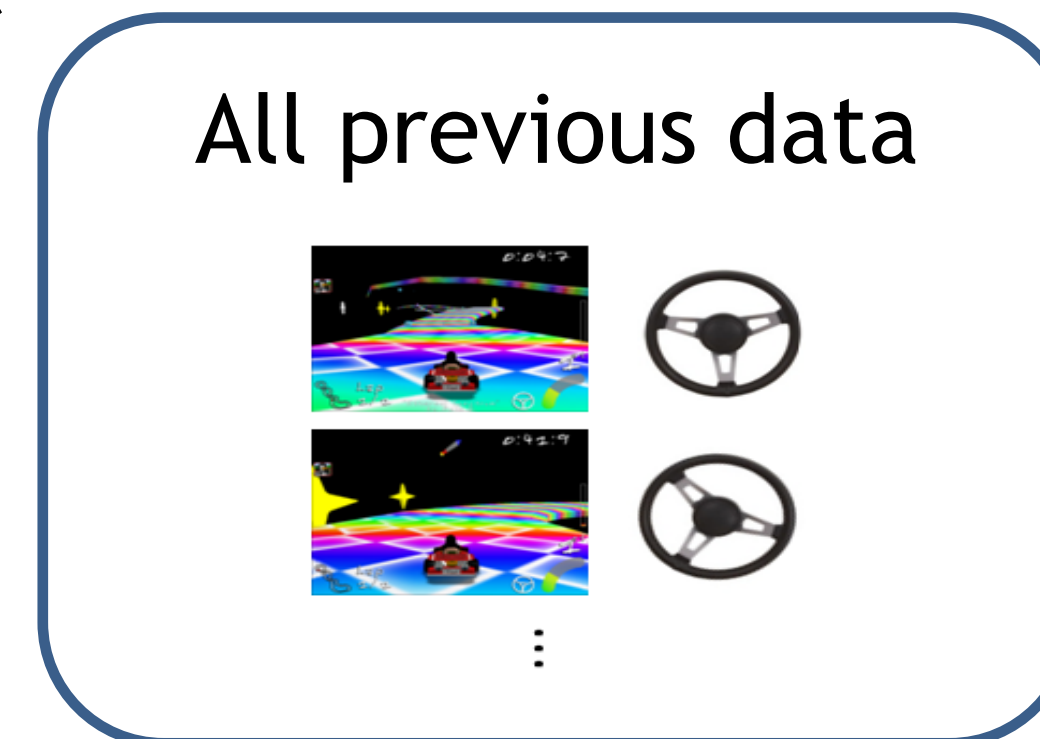
Steering
from
expert



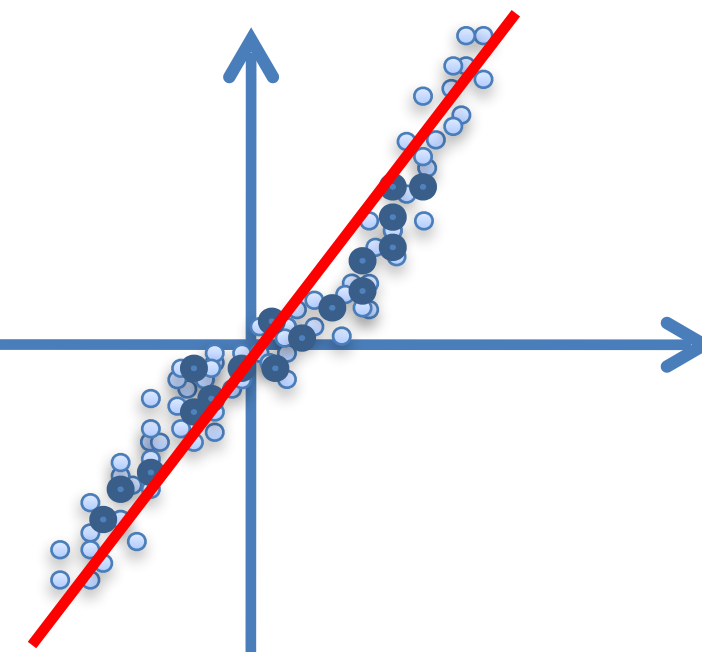
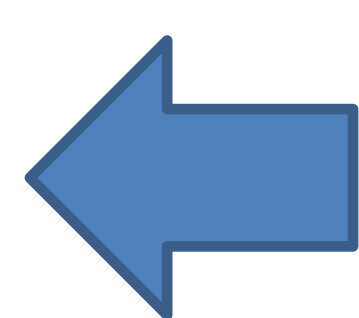
New Data



Aggregate
Dataset



New policy
 π_n



Supervised Learning

The DAgger algorithm

Initialize π^0 , and dataset $\mathcal{D} = \emptyset$

For $t = 0 \rightarrow T - 1$:

1. W/ π^t , generate dataset of trajectories $\mathcal{D}^t = \{\tau_1, \tau_2, \dots\}$

where for all trajectories $s_h \sim \rho_{\pi^t}$, $a_h = \pi^*(s_h)$

2. **Data aggregation:** $\mathcal{D} = \mathcal{D} \cup \mathcal{D}^t$

3. **Update policy via Supervised-Learning:** $\pi^{t+1} = \text{SL}(\mathcal{D})$

In practice, the DAgger algorithm requires less human labeled data than BC.

[\[Informal Theorem\]](#) Under more assumptions + assuming ϵ SL error is achievable, the DAgger algorithm has error: $|V^{\pi^*} - V^{\hat{\pi}}| \leq H\epsilon$

Success!

[Ross AISTATS 2011]



Today

- ✓ • Feedback from last lecture
- ✓ • Recap
- ✓ • Imitation Learning problem statement
- ✓ • Behavioral Cloning
- ✓ • DAgger

Summary:

1. IL can help a lot to explore the space
2. BC pretty good but brittle -> quadratic-in-horizon error
3. Online expert feedback can help with robustness -> linear-in-horizon error

Attendance:

bit.ly/3RcTC9T



Feedback:

bit.ly/3RHtlxy

