# Imitation Learning & Behavioral Cloning

**Lucas Janson and Sham Kakade**

**CS/Stat 184: Introduction to Reinforcement Learning**
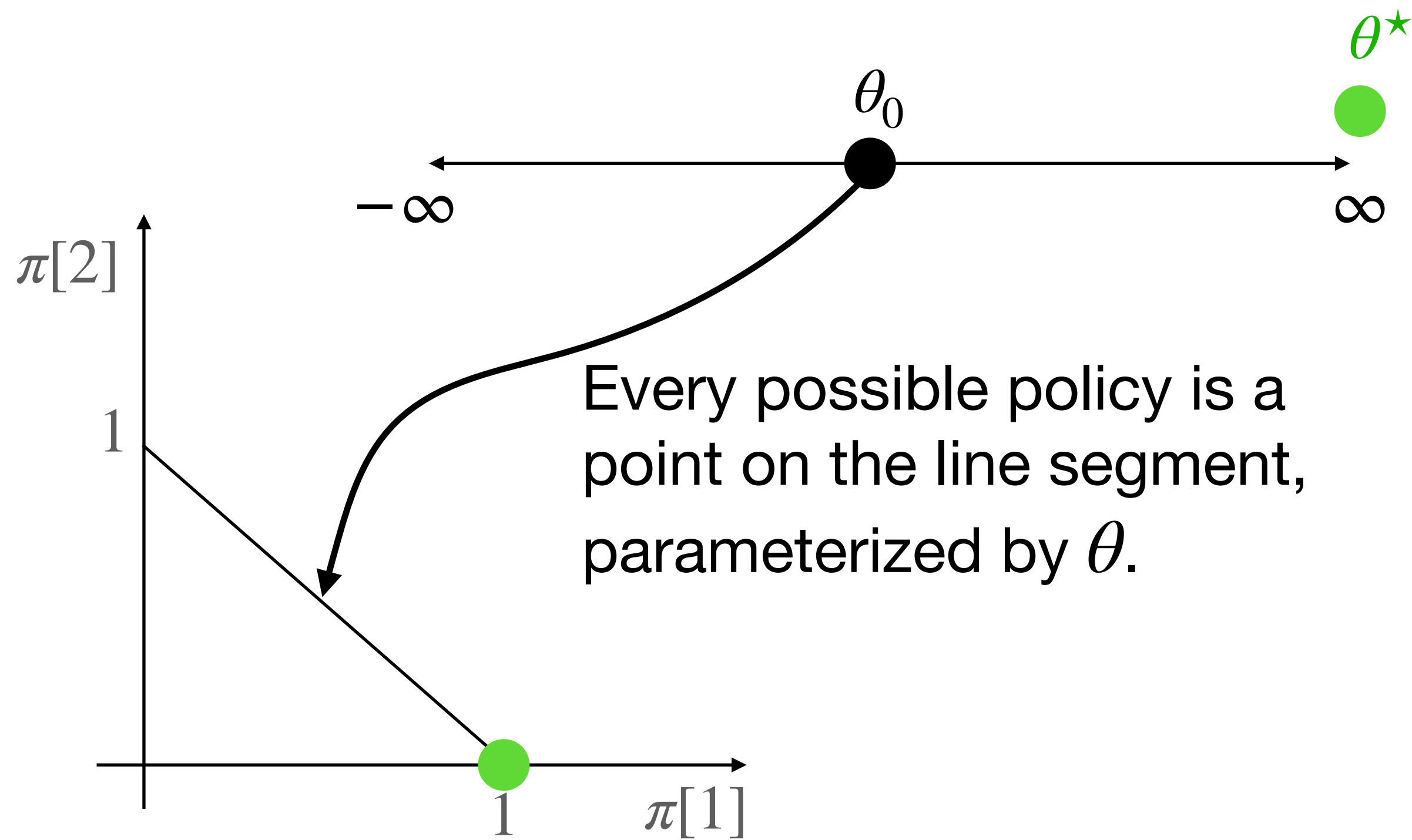**Fall 2023**

# Today

✓ • Recap++

• Imitation Learning:

    • Behavioral Cloning

    • DAgger

# Recap

# Example of Natural Gradient on 1-d problem: 2 actions, 1 state

$$(\pi_\theta[1], \pi_\theta[2]) := \left( \frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

$$J(\theta) = 100 \cdot \pi_\theta[1] + 1 \cdot \pi_\theta[2]$$

Gradient: $J'(\theta) = \dfrac{99 \exp(\theta)}{(1 + \exp(\theta))^2}$

Exact PG: $\theta^{k+1} = \theta^k + \eta \dfrac{99 \exp(\theta^k)}{(1 + \exp(\theta^k))^2}$

i.e., vanilla GA moves to $\theta = \infty$ with smaller and smaller steps, since $J'(\theta) \to 0$ as $\theta \to \infty$

Fisher information scalar: $F_\theta = \dfrac{\exp(\theta)}{(1 + \exp(\theta))^2}$

NPG: $\theta^{k+1} = \theta^k + \eta \dfrac{J'(\theta^k)}{F_{\theta^k}} = \theta_k + \eta \cdot 99$

NPG moves to $\theta = \infty$ much more quickly (for a fixed $\eta$)

$\theta^\star$

$\theta_0$

$-\infty$          $\infty$

$\pi[2]$

Every possible policy is a point on the line segment, parameterized by $\theta$.

1

$\pi[1]$

# Meta-Approach: CPI/TRPO/NPG/PPO are all pretty similar.

1. Init $\pi_0$

2. For $k = 0, \ldots K$:

$$\pi^{k+1} \approx \arg\max_{\theta} \Delta_k(\pi^\theta), \qquad \text{where } \Delta_k(\pi) = \mathbb{E}_{s_0, \ldots s_{H-1} \sim \rho_{\pi^k}} \left[ \sum_{h=0}^{H-1} \mathbb{E}_{a_h \sim \pi(s_h)} A^{\pi^k}(s_h, a_h) \right]$$

   such that $\rho_\theta$ is "close" to $\rho_{\theta^k}$

   - CPI: conservative policy iteration
     uses unconstrained optimization: $\widetilde{\pi} \approx \arg\max_{\theta} \Delta_k(\pi^\theta)$,

     enforces closeness with "mixing": $\pi^{k+1} = (1 - \alpha) \cdot \pi^k + \alpha \cdot \widetilde{\pi}^{k+1}$
   - TRPO: use KL to enforce closeness.
   - NPG: is TRPO up to "leading order" (via Taylor's theorem).
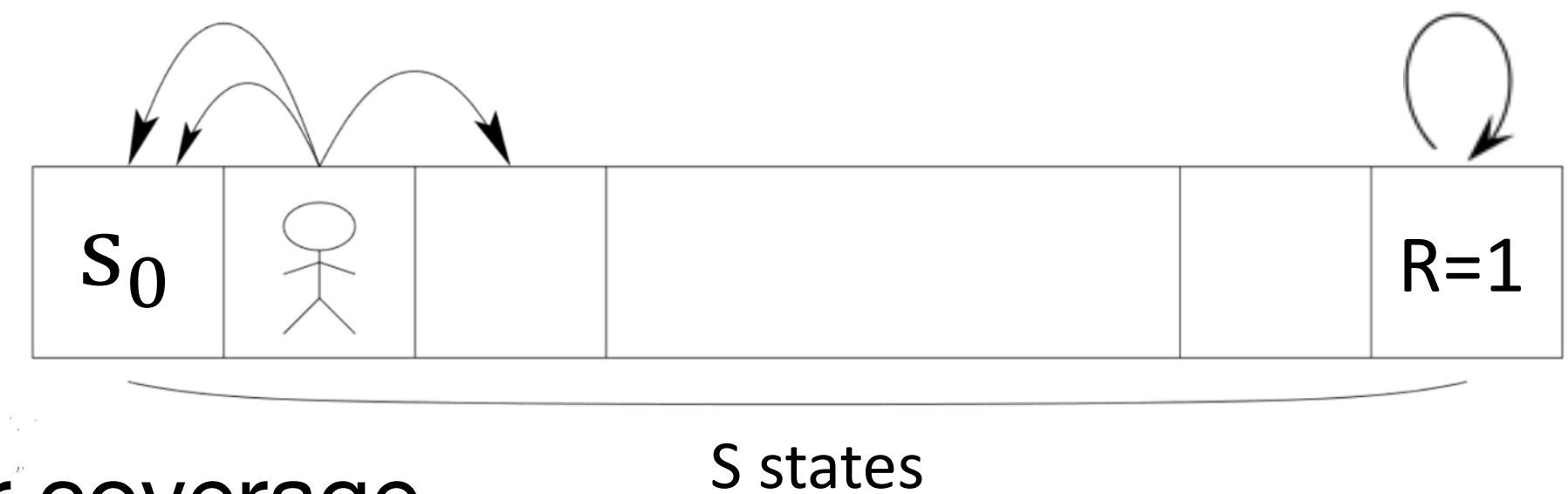   - PPO: uses a Lagrangian relaxation (i.e. regularization)

3. Return $\pi_K$

# "Lack of Exploration" leads to Optimization and Statistical Challenges



$s_0$ ... R=1

S states

Thrun '92

- Suppose $H \approx \text{poly}(|S|)$ & $\mu(s_0) = 1$ (i.e. we start at $s_0$).
- A randomly initialized policy $\pi^0$ has prob. $O(1/3^{|S|})$ of hitting the goal state in a trajectory.
- Implications:
  - The following sample based approach, with $\mu(s_0) = 1$, require $O(3^{|S|})$ trajectories.
    - Holds for (sample based) Fitted DP
    - Holds for (sample based) PG/CPI/TRPO/NPG/PPO
- Basically, for these approaches, we are stuck without exploration, if $\mu(s_0) = 1$.

# Let's examine the role of $\mu$



S states

Thrun '92

- Suppose that somehow the distribution $\mu$ had better coverage.

  - e.g, $\mu$ was uniform over the all states in our toy problem, then all approaches we covered would work (with mild assumptions )
  - Theory: CPI/TRPO/NPG/PPO have better guarantees than fitted DP methods (assuming some "coverage")
- Strategies:
  - If we have a simulator, sometimes we can design $\mu$ to have better coverage.
    - this is helpful for robustness as well.
  - Imitation learning

    - An expert gives us samples from a "good" $\mu$.
  - Explicit exploration:
    - UCB-VI: we'll merge two good ideas!
    - Encourage exploration in PG methods.
  - Try with reward shaping

# Today:

# What about guarantees for PG methods? (vs fitted-DP methods)

- The hope is that if (average case) "supervised learning" worked, then RL would also work.

$$\pi_\theta(a \,|\, s) = \frac{\exp(\theta^\top \phi(s, a))}{\sum_{a'} \exp(\theta^\top \phi(s, a'))}$$
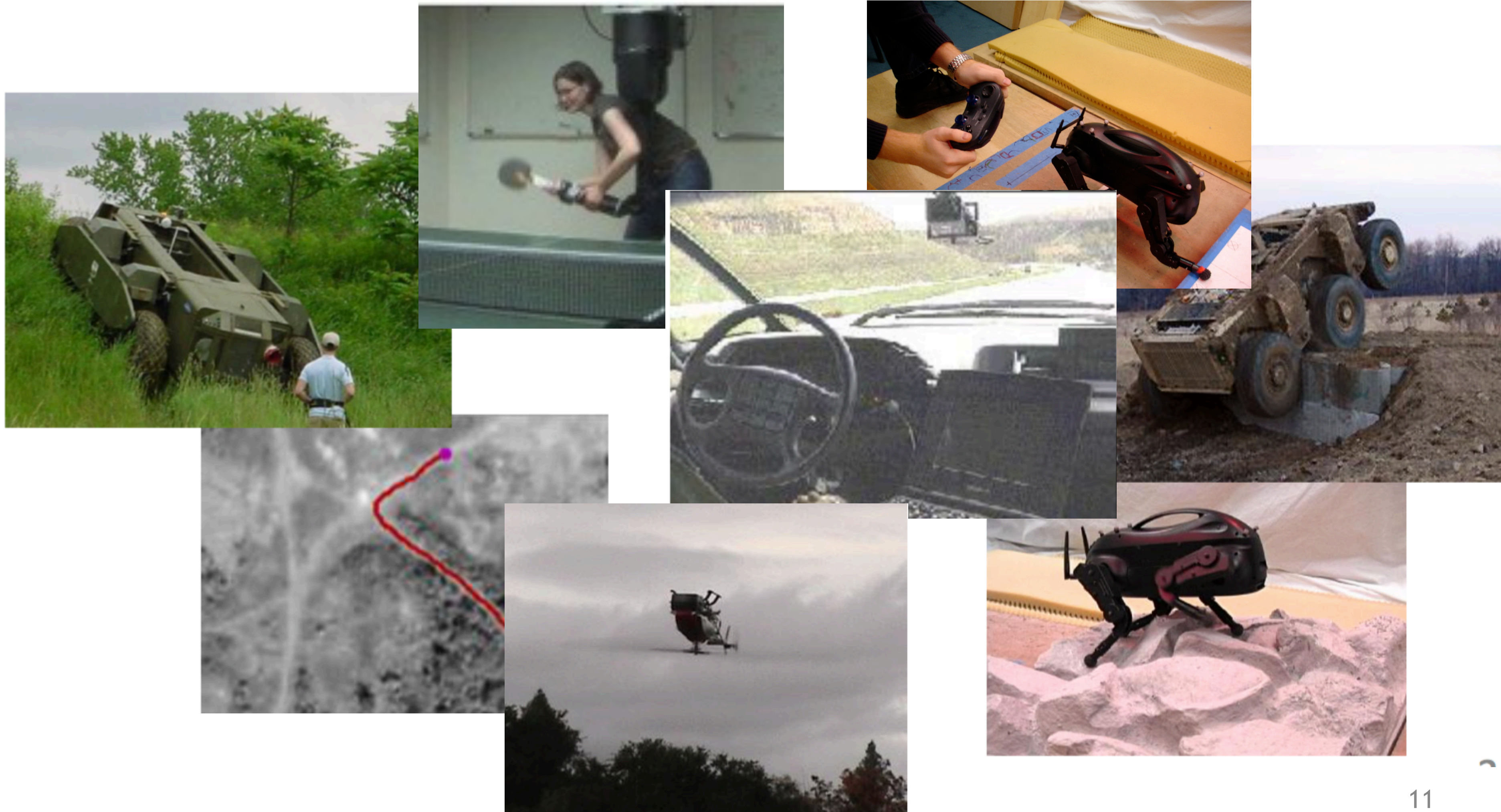
Issues: let's consider log-linear policies.

- Approximation error: For log linear policies, how good does $\phi$ need to be?
  (comment: hopefully some average case condition for approximating $A^\pi(s, a)$)
- Sample size: hope to use a # samples that is poly in dim($\phi$) $\&$ $1/\epsilon_{accuracy}$.
- Coverage: need some coverage condition over the state space.
  (comment: hopefully the coverage conditions are only in "$\phi$-space")
- Computation: we want NPG to find something good with poly in $d, 1/\epsilon_{accuracy}, H$ iterations.

- **Theory**: (see AJKS Ch 4+13, for formal log linear policies)
  There are (somewhat subtle) approx/coverage conditions where NPG converges to an $\epsilon_{accuracy}$-opt policy with
  poly sample, poly computation time.
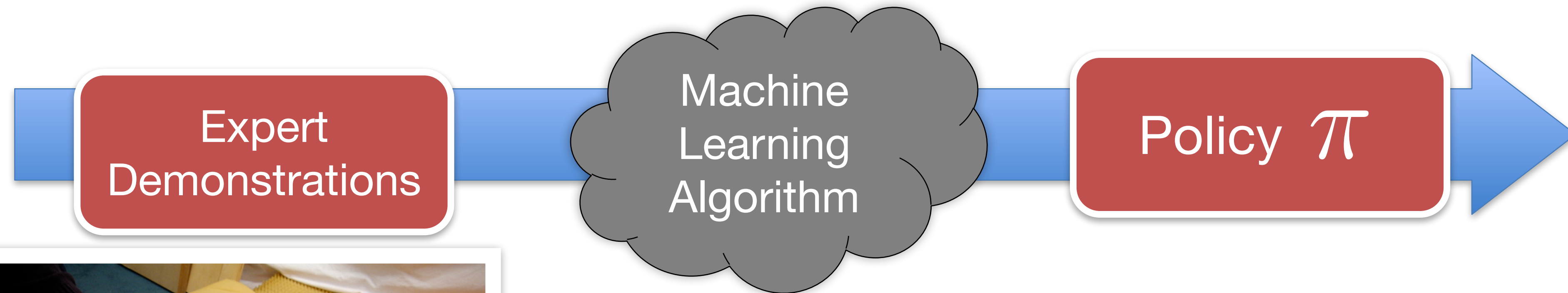  (Conditions are weaker than those for fitted-DP methods).

# Today

- Recap++

✔ - Imitation Learning:

  - Behavioral Cloning

  - DAgger

# Imitation Learning

# Imitation Learning



**Expert Demonstrations** → Machine Learning Algorithm → **Policy $\pi$** →

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

Maps *states* to <u>actions</u>

# Learning to Drive by Imitation

[Pomerleau89, Saxena05, Ross11a]

## Input:



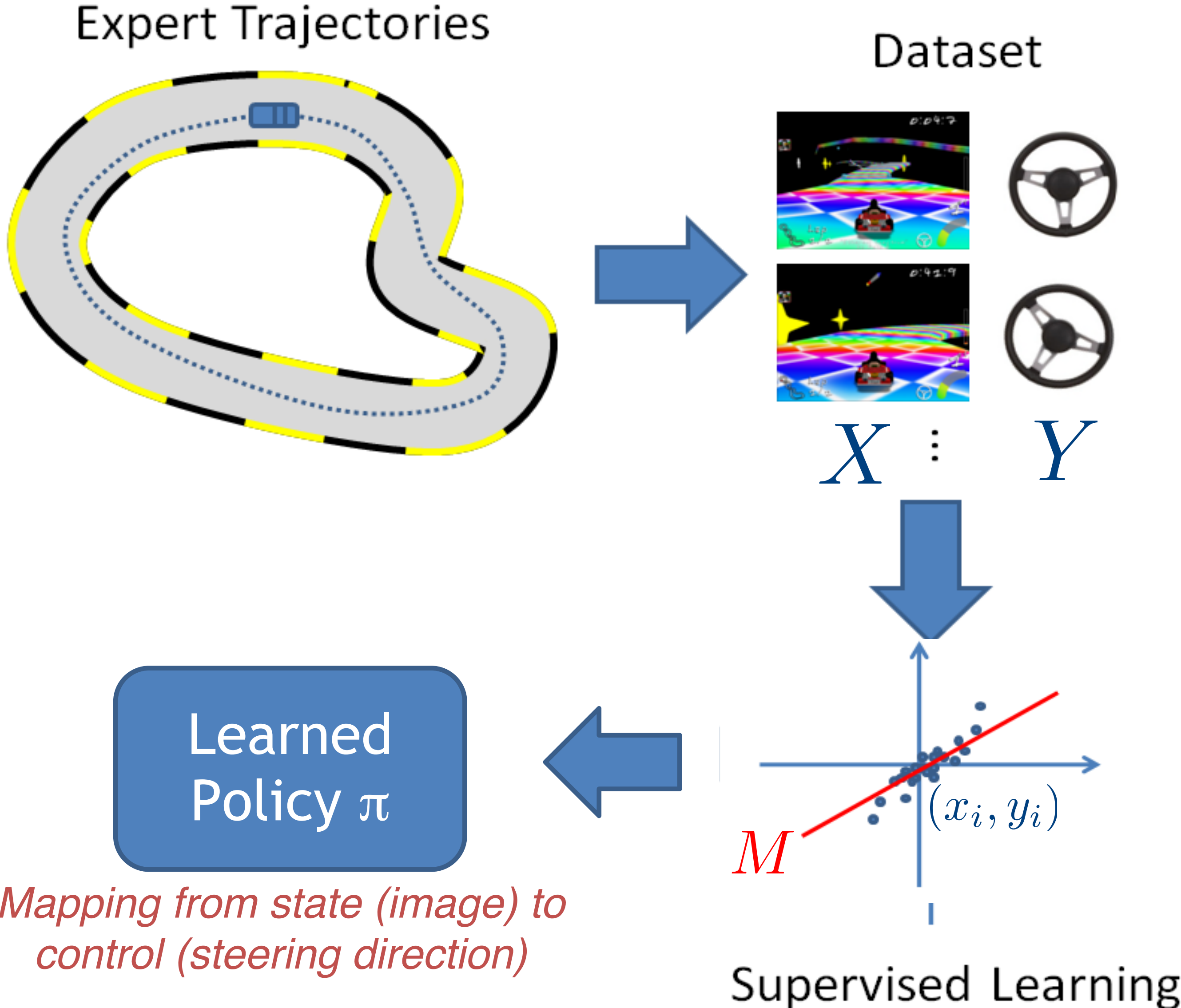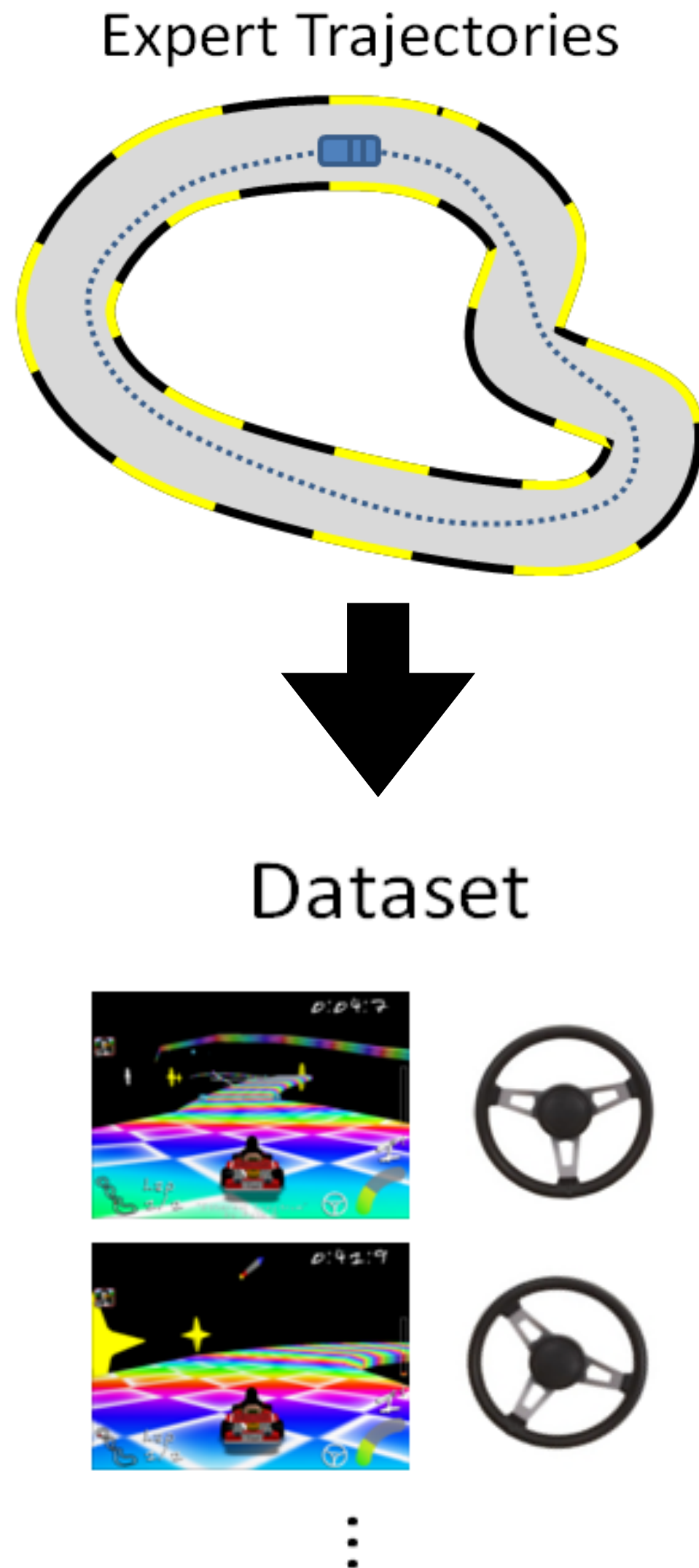Camera Image

## Output:

Policy



Steering Angle
in [-1, 1]

# Today

- Recap++

- Imitation Learning:

  ✓ - Behavioral Cloning

  - DAgger

# Supervised Learning Approach: Behavior Cloning

Expert Trajectories

Dataset

$X \quad\vdots\quad Y$



**Learned Policy π**

*Mapping from state (image) to control (steering direction)*

$M$    $(x_i, y_i)$

Supervised Learning

15

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

Expert Trajectories

Dataset

Finite horizon MDP $\mathcal{M}$

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

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

Goal: learn a policy from $\mathcal{D}$ that is as good as the expert $\pi^\star$

# Let's formalize the Behavior Cloning algorithm

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

BC is a Reduction to Supervised Learning:

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

Many choices of loss functions:

1. Negative log-likelihood (NLL): $\ell(\pi, s, a) = -\ln \pi(a \mid s)$

2. 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

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

Note a training and testing "mismatch"

Theorem [BC Performance]:

suppose we assume supervised learning succeeds, with $\epsilon$ classification error:

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

(where $\pi^\star$ is the experts policy, which need not be optimal)

then, under $\mu$, we have:

$$|V^{\pi^\star} - V^{\widehat{\pi}}| \leq H^2 \epsilon$$

The quadratic amplification is annoying

18

# Proof:

By the PDL

$$|V^{\pi^\star}(s) - V^{\widehat{\pi}}(s)| = \left| \mathbb{E}_{\tau \sim \rho_{\pi^\star}} \left[ \sum_{h=0}^{H-1} A_h^{\widehat{\pi}}(s_h, a_h) \right] \right|$$

$$= \left| \mathbb{E}_{s_1, \ldots s_h \sim \rho_{\pi^\star}} \left[ \sum_{h=0}^{H-1} A_h^{\widehat{\pi}}(s_h, \pi^\star(s_h)) \right] \right|$$

$$\leq H \left| \mathbb{E}_{\tau \sim \rho_{\pi_{\pi^\star}}} \left[ \sum_{h=0}^{H-1} \mathbf{1} \left[ \widehat{\pi}(s_h) \neq \pi^\star(s_h) \right] \right] \right|$$

$$\leq H^2 \epsilon$$

# **Distribution Shift Example** ($H^2$ factor is tight)



$r(s_1) = 1$

Initial state

Opt policy:

Under $\rho_{\pi^\star}$, trajectory is $s_0, s_1, s_1, \ldots$

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

$V_H^{\pi^\star}(s_0) = H - 1$

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

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

This policy has good supervised learning error:

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

note: while $\widehat{\pi}(s_2) \neq \widehat{\pi}^\star(s_2)$, state $s_2$ is never visited under $\pi^\star$

We have quadratic degradation (in $H$):

$$V_H^{\pi^\star}(s_0) = (1 - H\epsilon) \cdot V_H^{\pi^\star}(s_0) + H\epsilon \cdot 0 = V_H^{\pi^\star}(s_0) - \epsilon H(H-1)$$

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

# What could go wrong?

- Predictions affect future inputs/ observations



Learned Policy

Expert's trajectory

# Expert Demos

# BC Policy

# Today

- Recap++

- Imitation Learning:

  - Behavioral Cloning

✓ - DAgger

# Intuitive solution: **Interaction**

**Use interaction to collect data where learned policy goes**

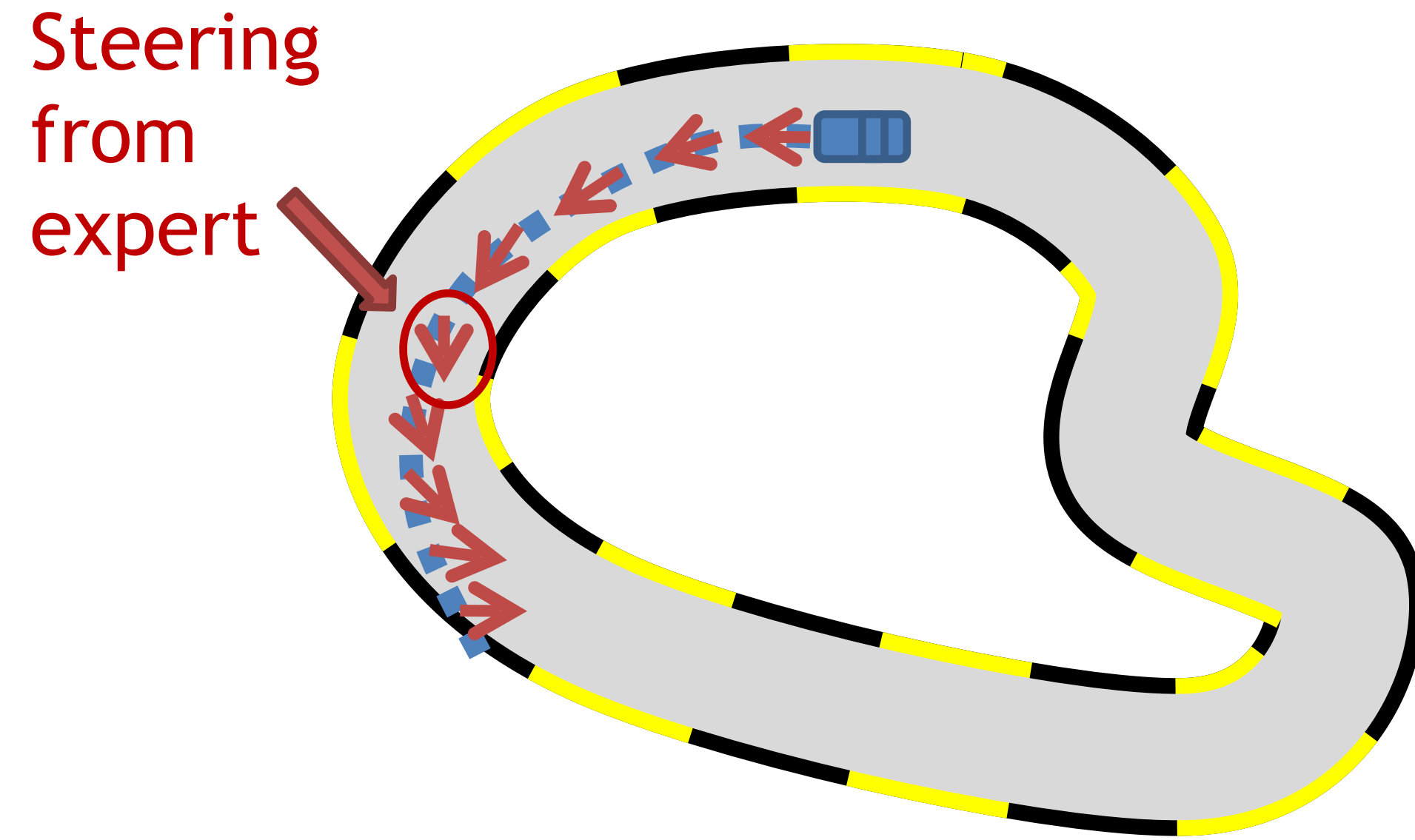# General Idea: Iterative Interactive Approach

New Data

Collect Data through Interaction

Update Policy

Updated Policy

# DAgger: Dataset Aggregation [Ross11a]
## 0th iteration

**Expert Demonstrates Task**

**Dataset**

**Supervised Learning**

1st policy $\pi_1$

# DAgger: Dataset Aggregation [Ross11a]

## 1st iteration

**Execute $\pi_1$ and Query Expert**

Steering from expert

# DAgger: Dataset Aggregation <superscript>[Ross11a]</superscript>

## 1st iteration

**Execute $\pi_1$ and Query Expert**

**New Data**

Steering from expert

States from
the learned policy

# DAgger: Dataset Aggregation [Ross11a]

## 1st iteration

**Execute $\pi_1$ and Query Expert**



Steering from expert

New Data

All previous data

# DAgger: Dataset Aggregation [Ross11a]
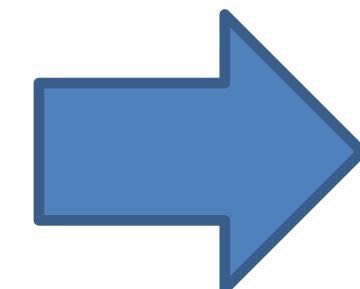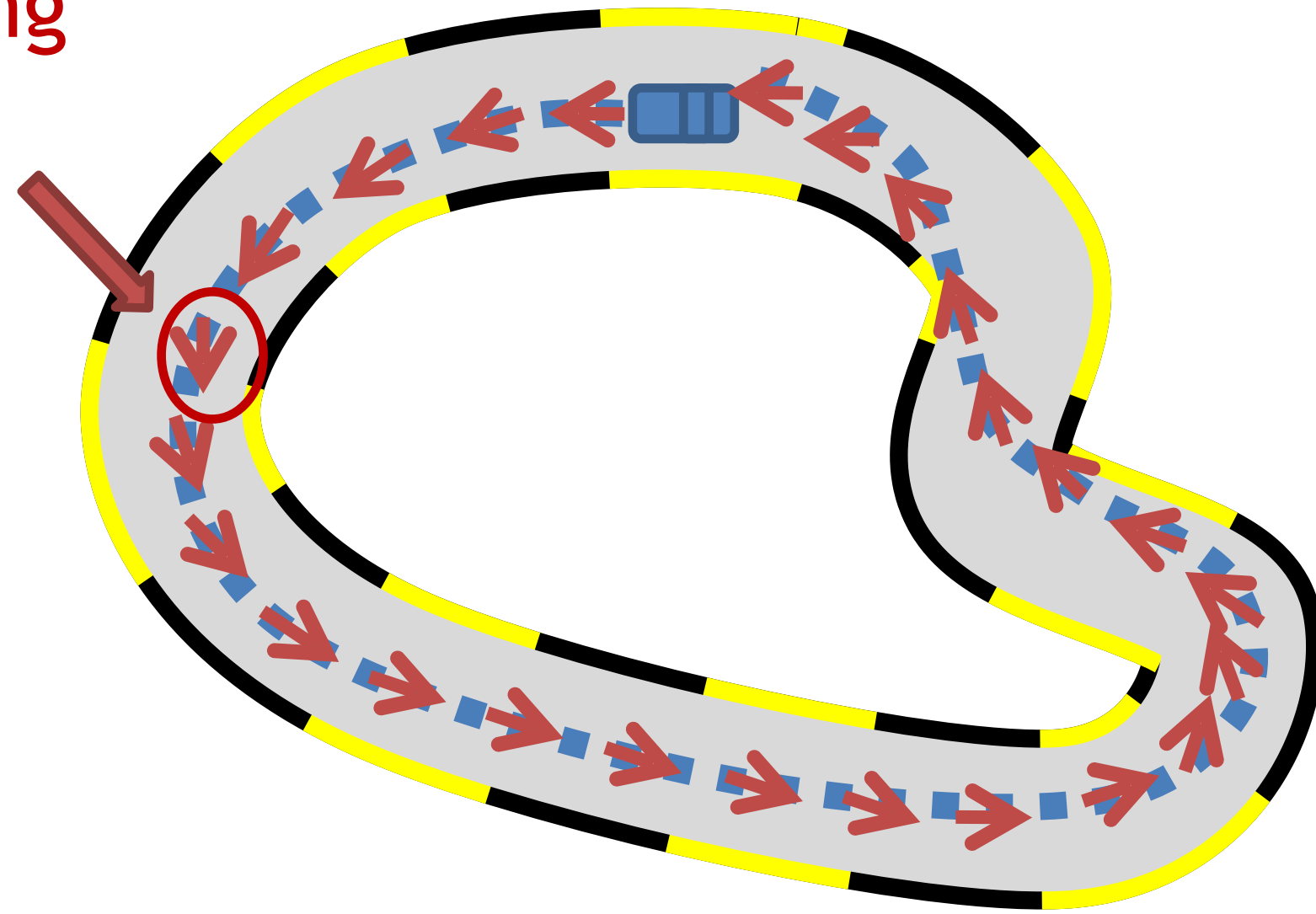## 1st iteration

**Execute $\pi_1$ and Query Expert**

Steering from expert

**New Data**

Aggregate Dataset

**+**

All previous data

New policy $\pi_2$

**Supervised Learning**

# DAgger: Dataset Aggregation [Ross11a]
## 2nd iteration

**Execute $\pi_2$ and Query Expert**

Steering from expert

**New Data**

Aggregate Dataset

All previous data

New policy $\pi_3$

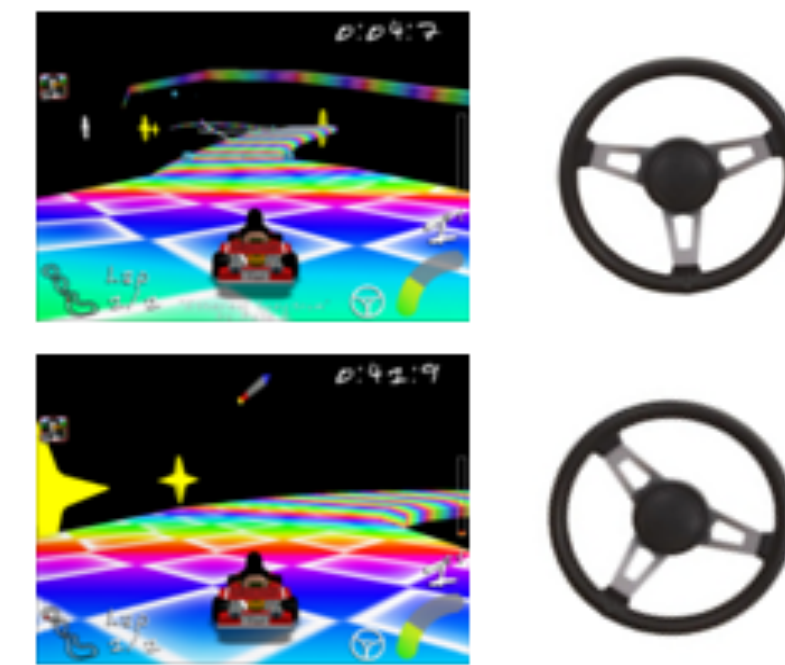**Supervised Learning**

# DAgger: Dataset Aggregation [Ross11a]

## nth iteration

**Execute $\pi_{n-1}$ and Query Expert**
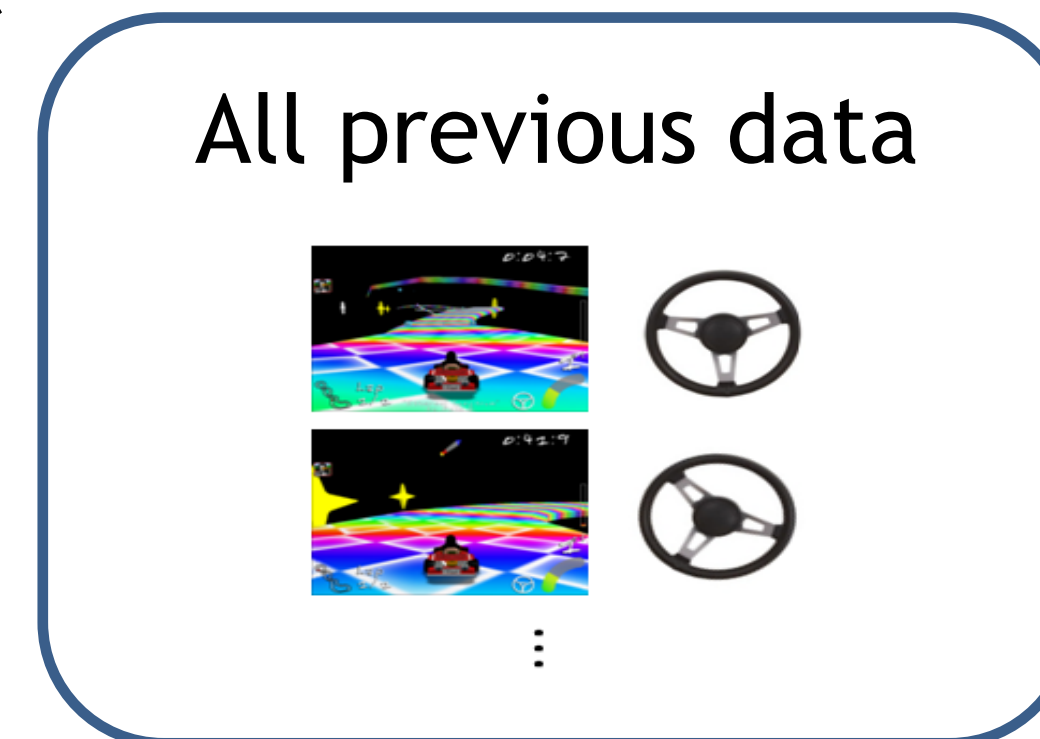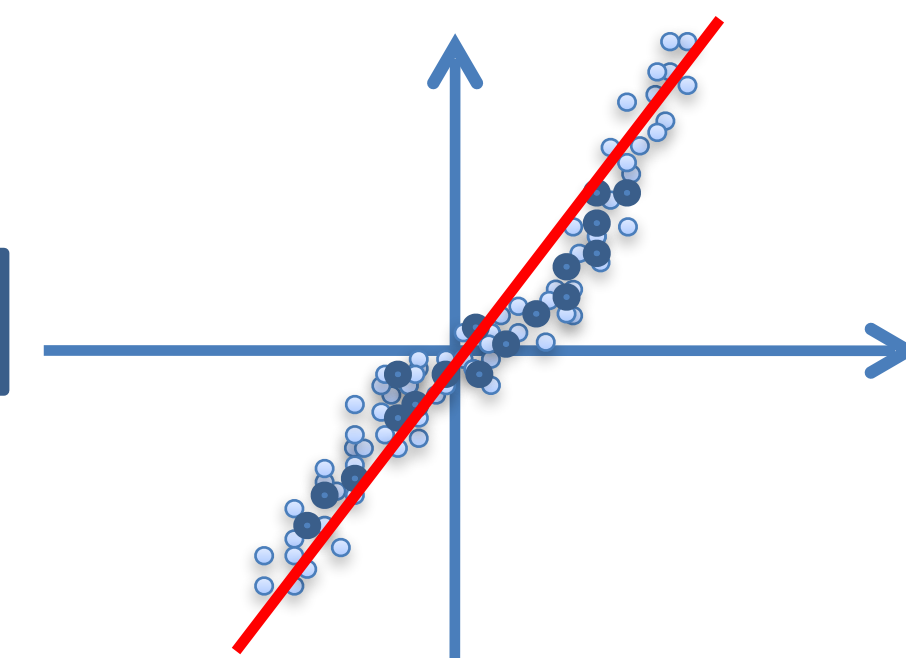
Steering from expert

**New Data**

Aggregate Dataset

All previous data

New policy $\pi_n$

**Supervised Learning**

# The DAgger algorithm

Initialize $\pi^0$, and dataset $\mathcal{D} = \varnothing$

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

    1. W/ $\pi^t$, generate dataset of trajectories $\mathcal{D}^t = \{\tau_1, \tau_2, \dots\}$ where for all trajectories $s_h \sim \rho_{\pi^t}$, $a_h = \textcolor{red}{\pi^\star(s_h)}$

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

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

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

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

# Success!

# Summary:

1. NPG: a simpler way to do TRPO, a "pre-conditioned" gradient method.
2. PPO: "first order" approx to TRPO

Attendance:
bit.ly/3RcTC9T

Feedback:
bit.ly/3RHtlxy