

Imitation Learning: Behavior Cloning, Distribution Shift, & Distribution Matching

Sham Kakade and Wen Sun

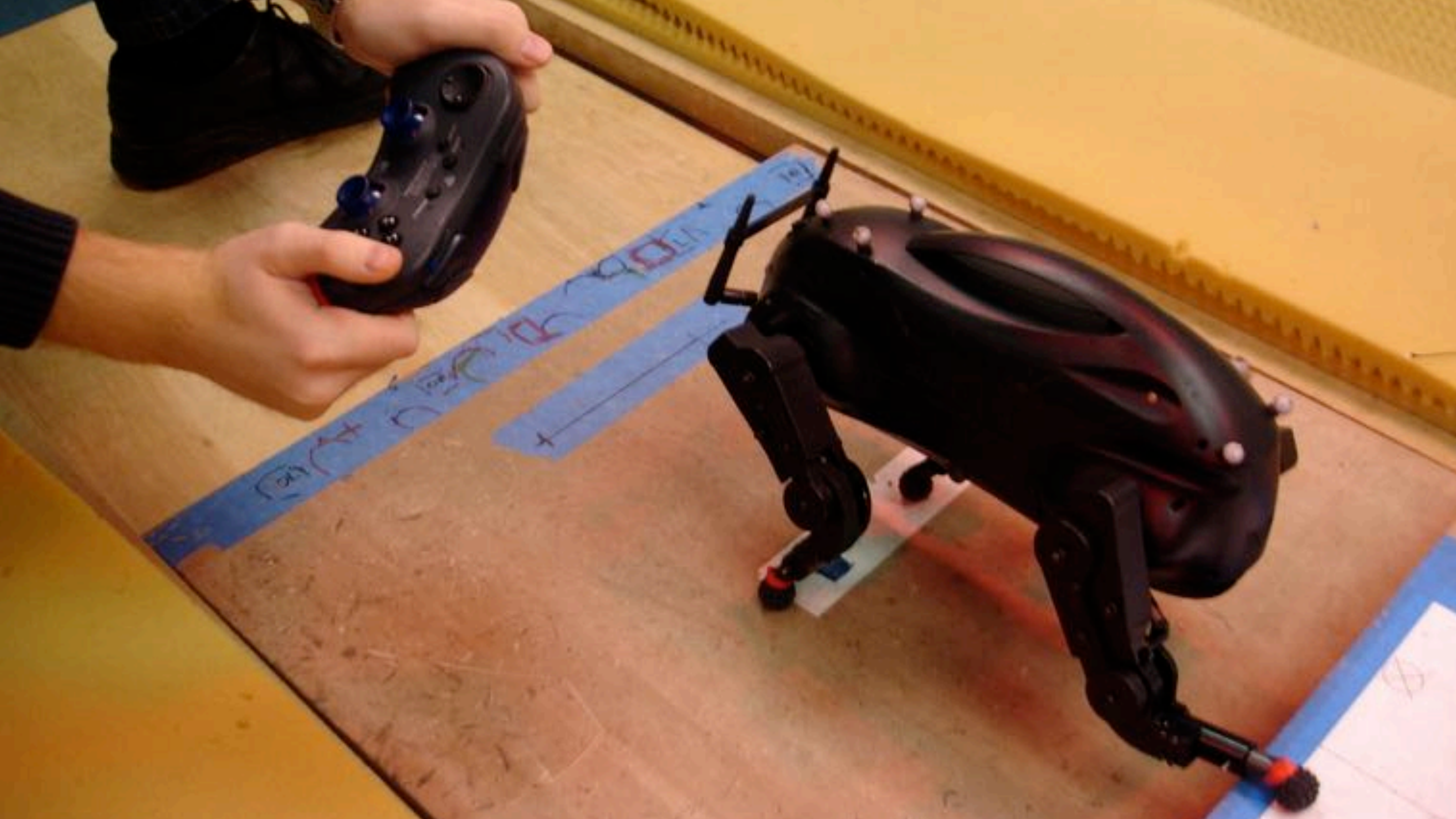
CS 6789: Foundations of Reinforcement Learning

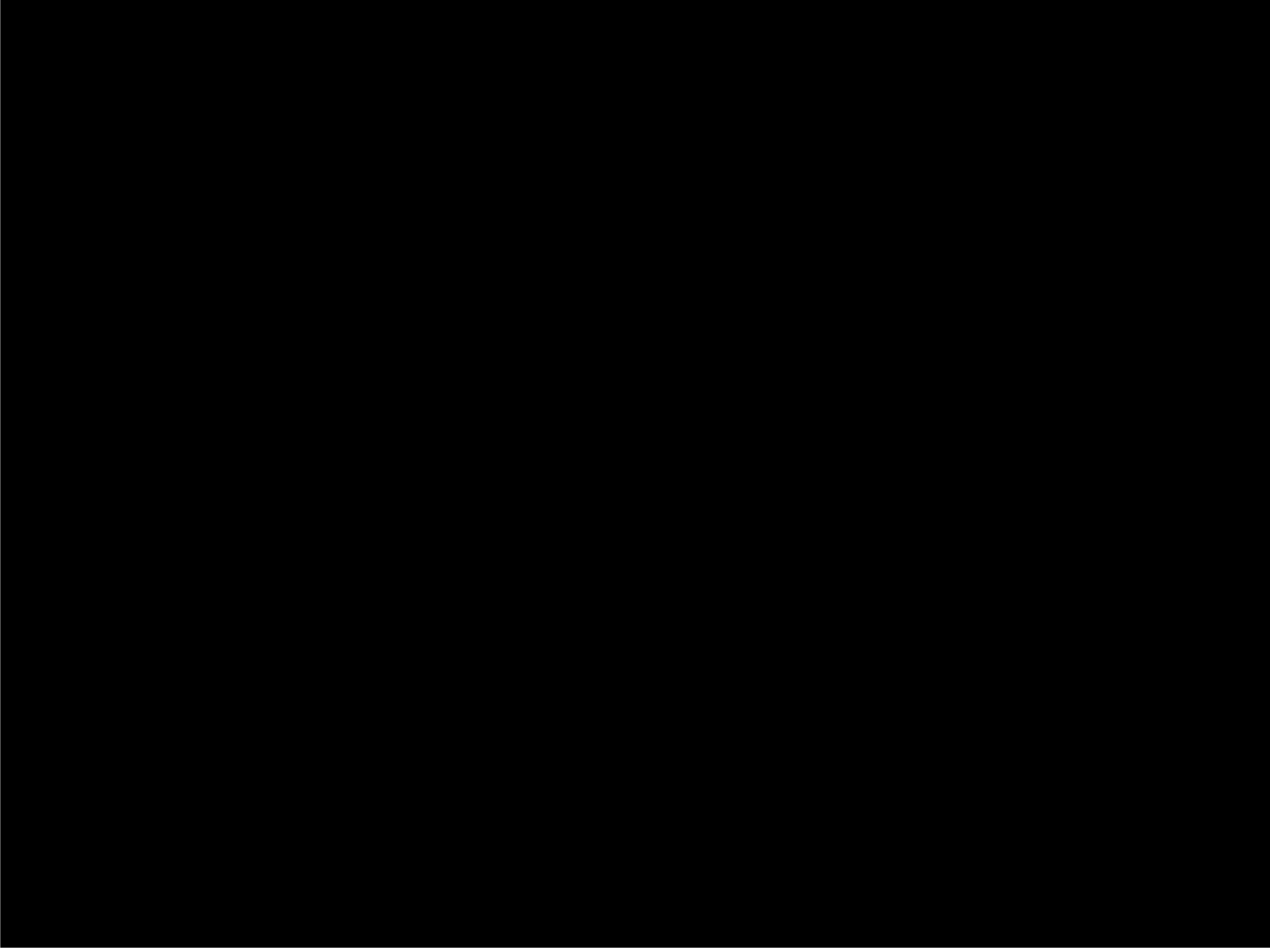
Today: Imitation Learning

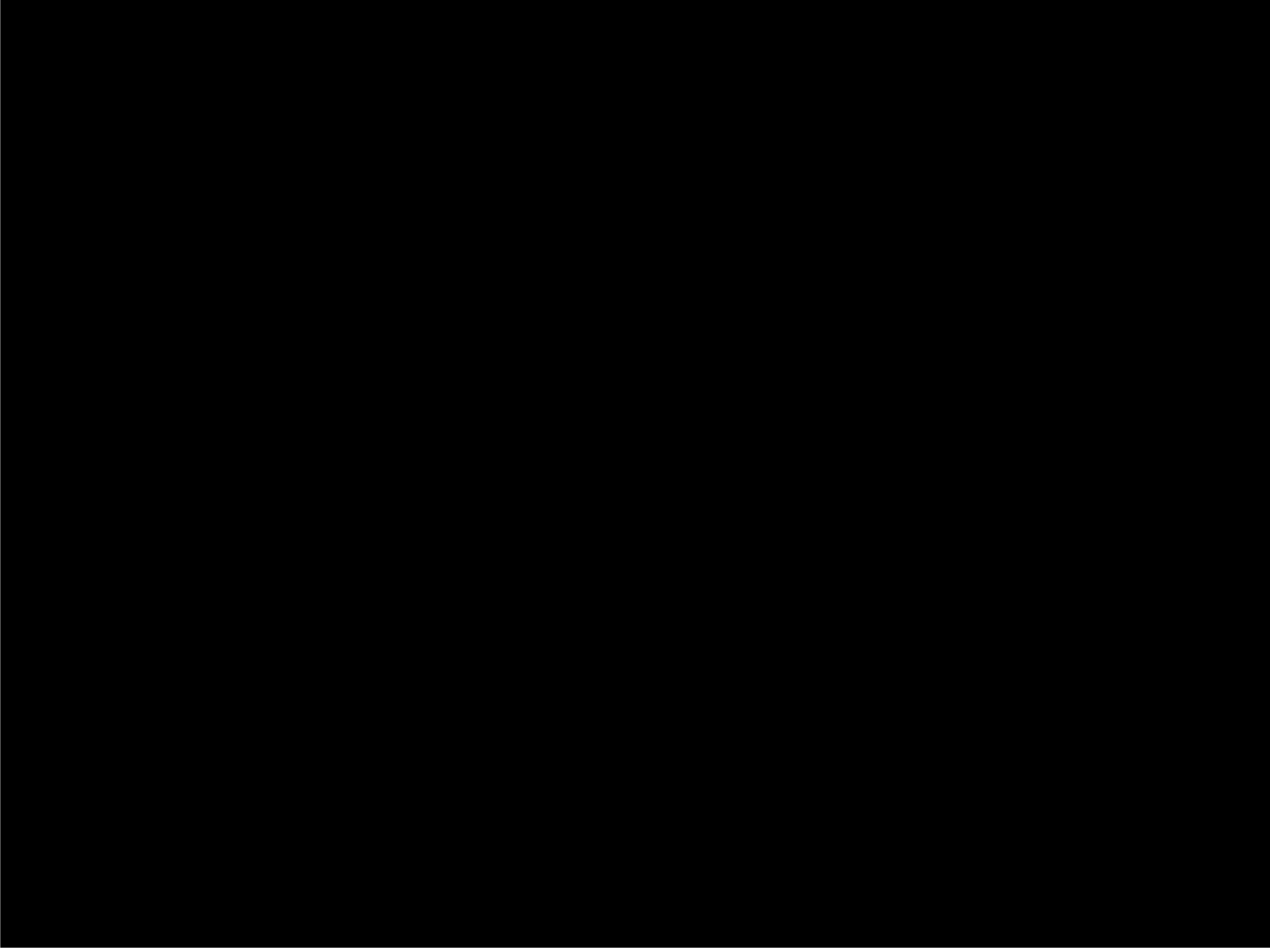
1. Introduction of Imitation Learning

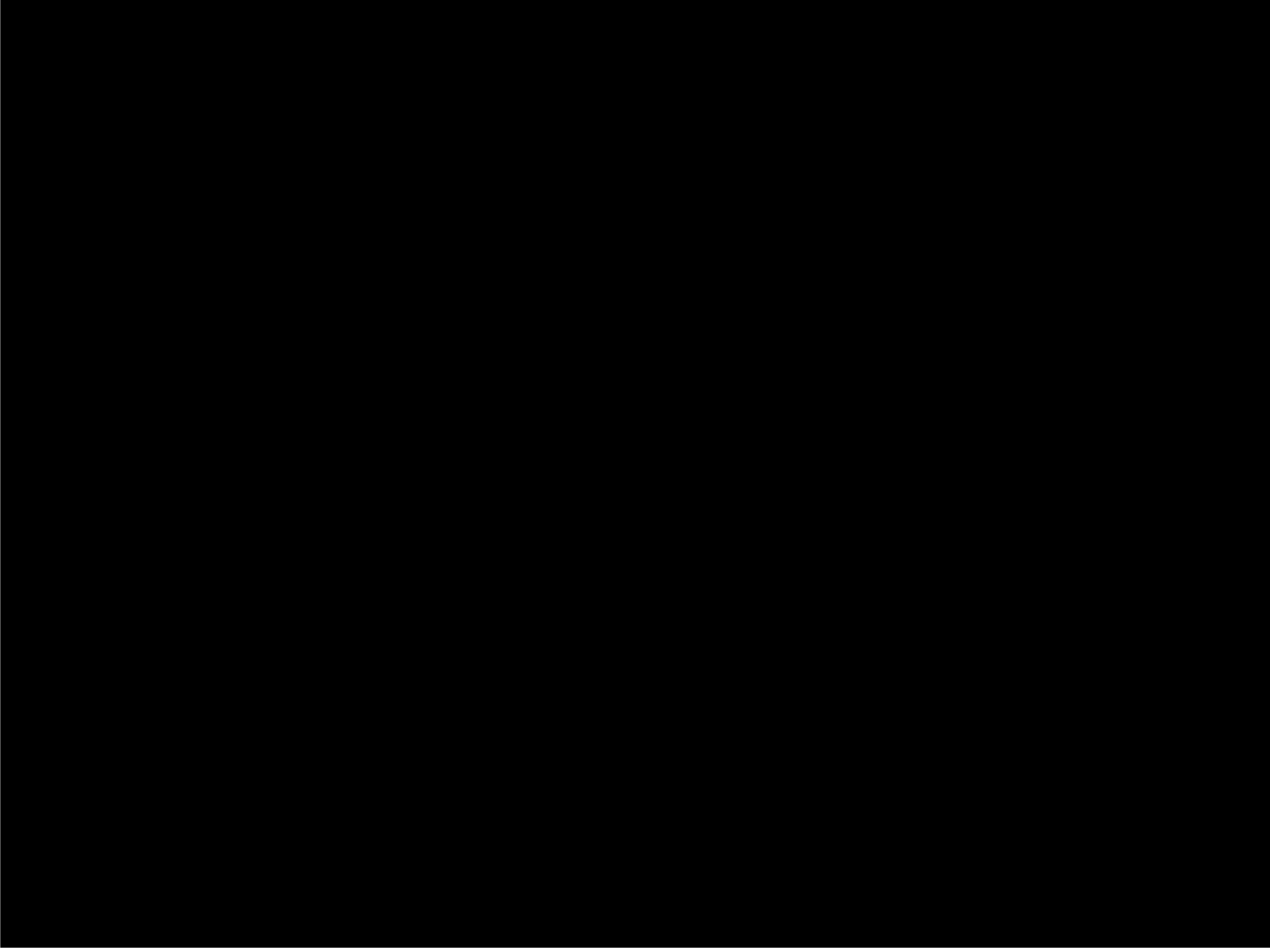
2. Offline Imitation Learning: Behavior Cloning

3. The hybrid Setting: Statistical Benefit and Distribution Matching









An Autonomous Land Vehicle In A Neural Network *[Pomerleau, NIPS '88]*

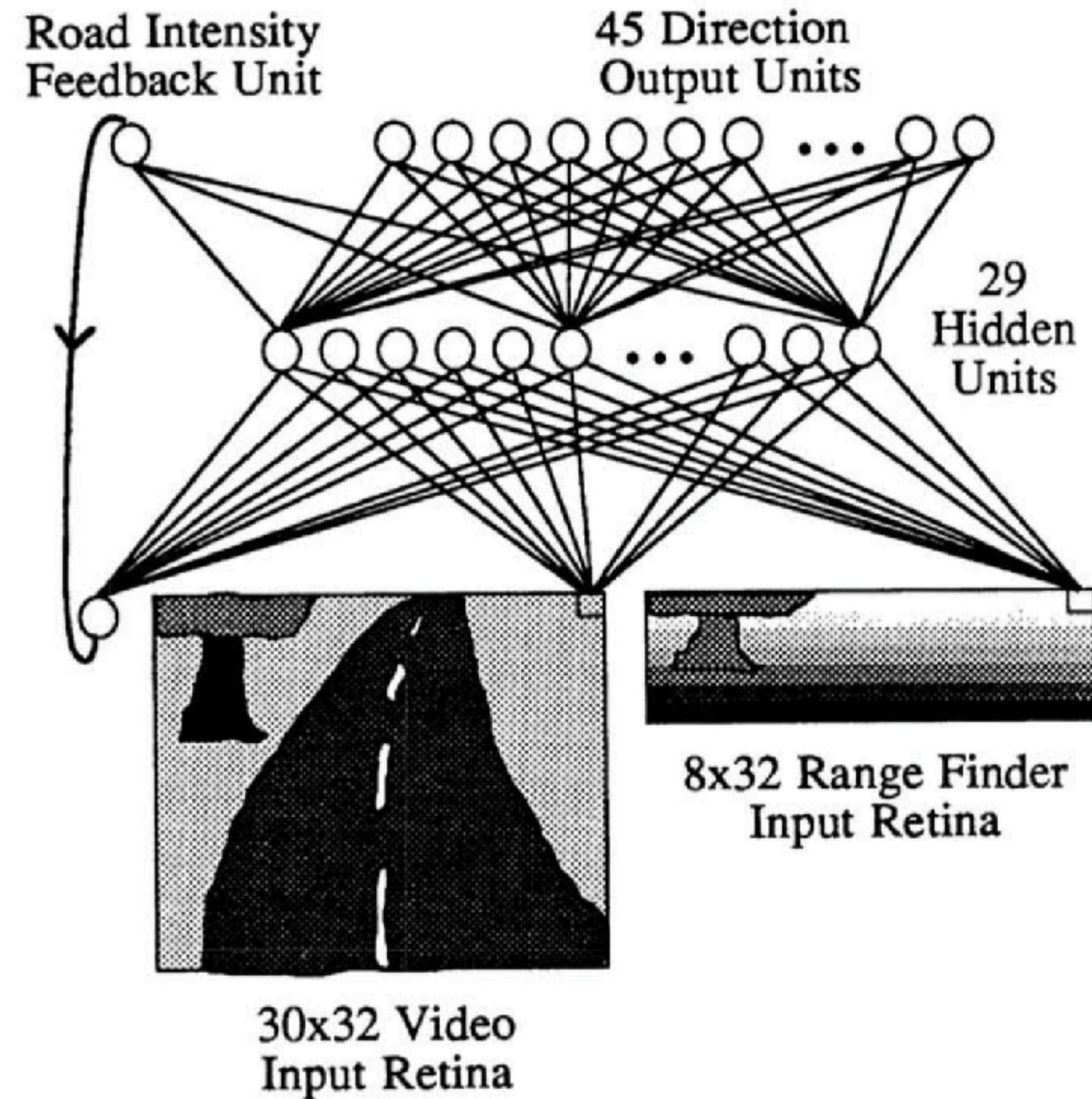
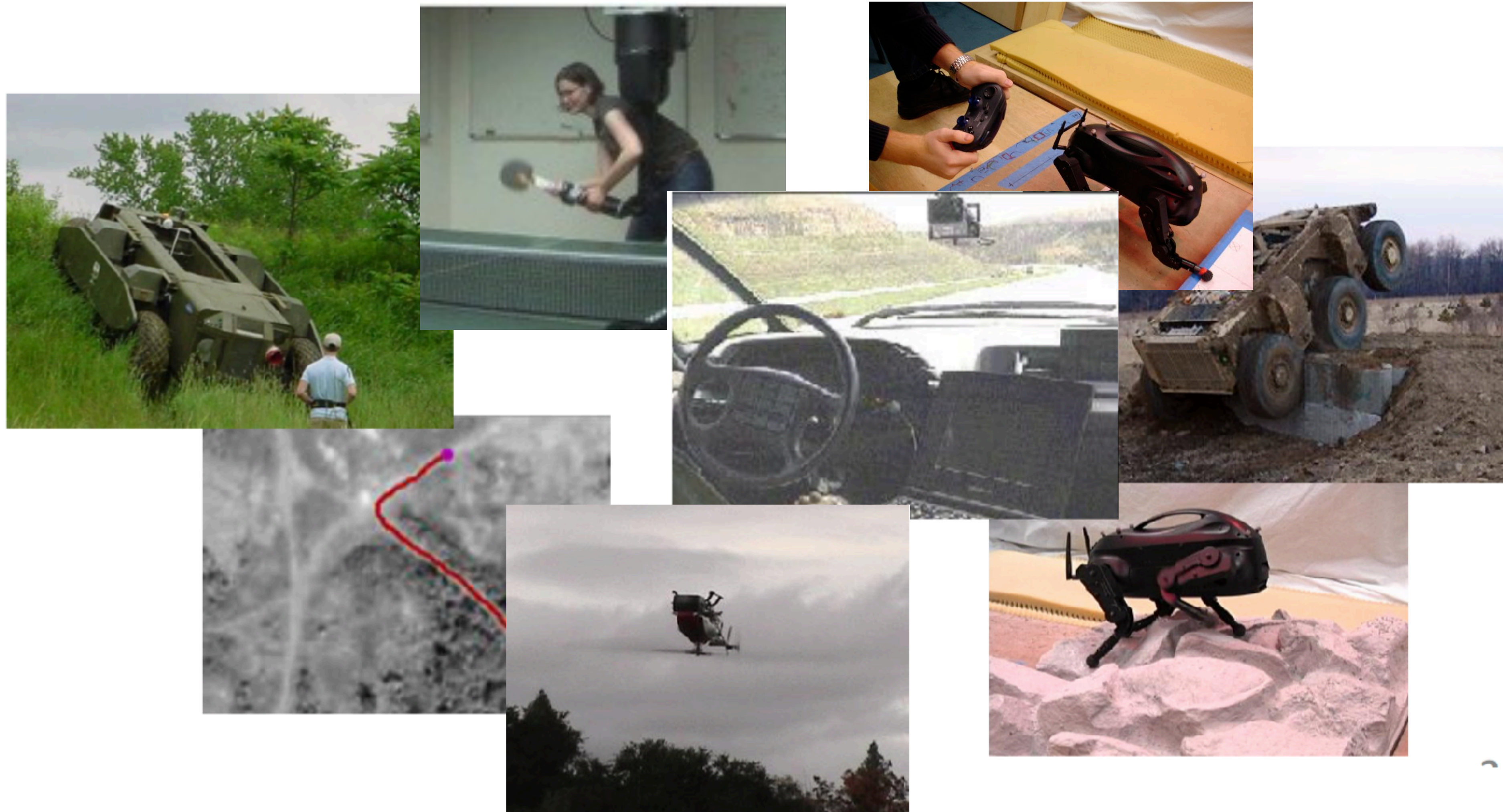


Figure 1: ALVINN Architecture

Imitation Learning



Imitation Learning

Imitation Learning



Imitation Learning

Expert
Demonstrations



Imitation Learning

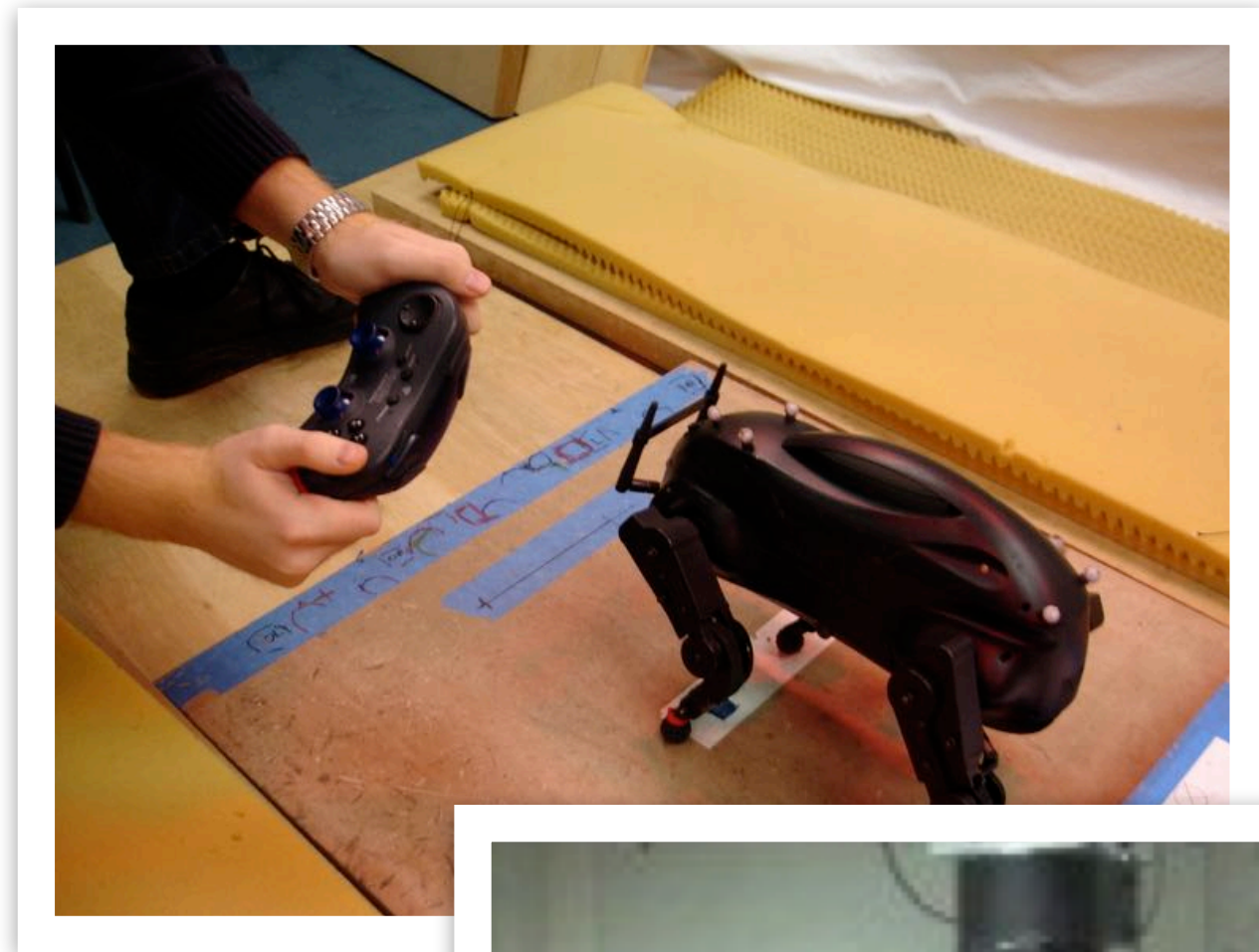
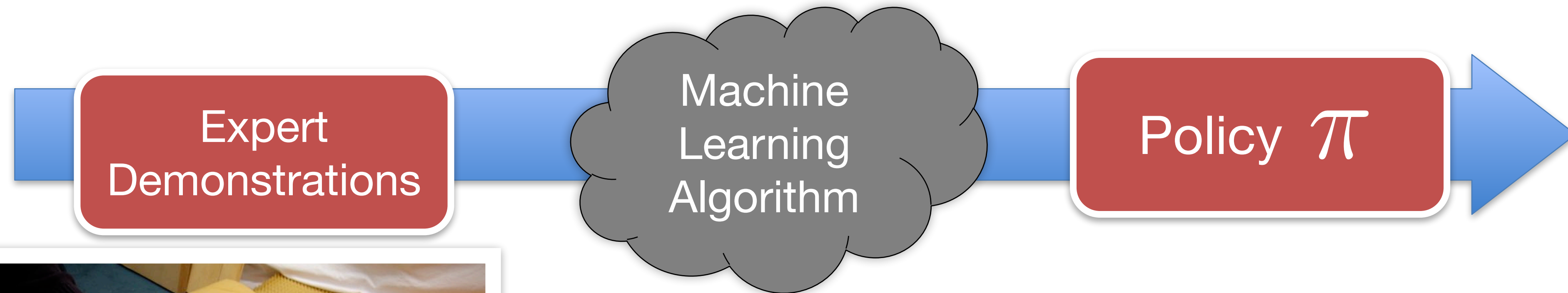
Expert
Demonstrations

Machine
Learning
Algorithm



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

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

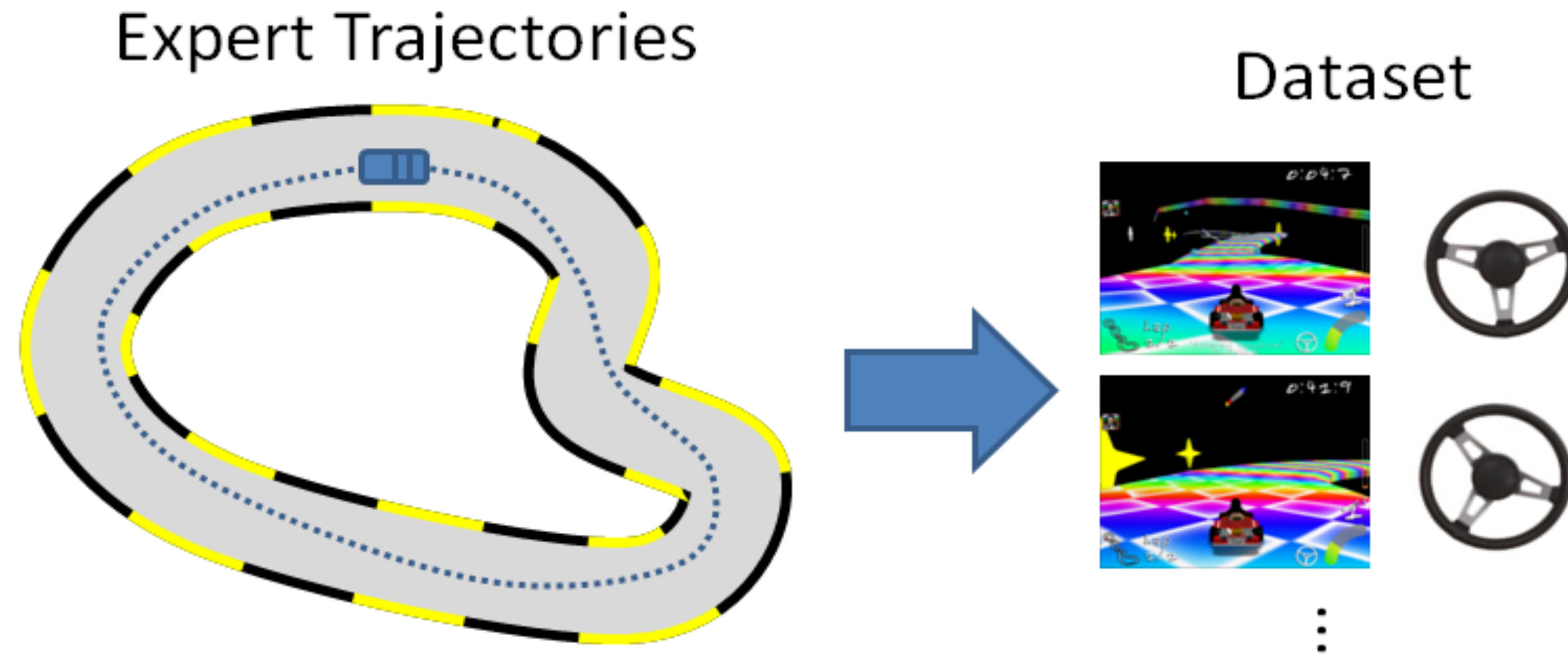


Output:

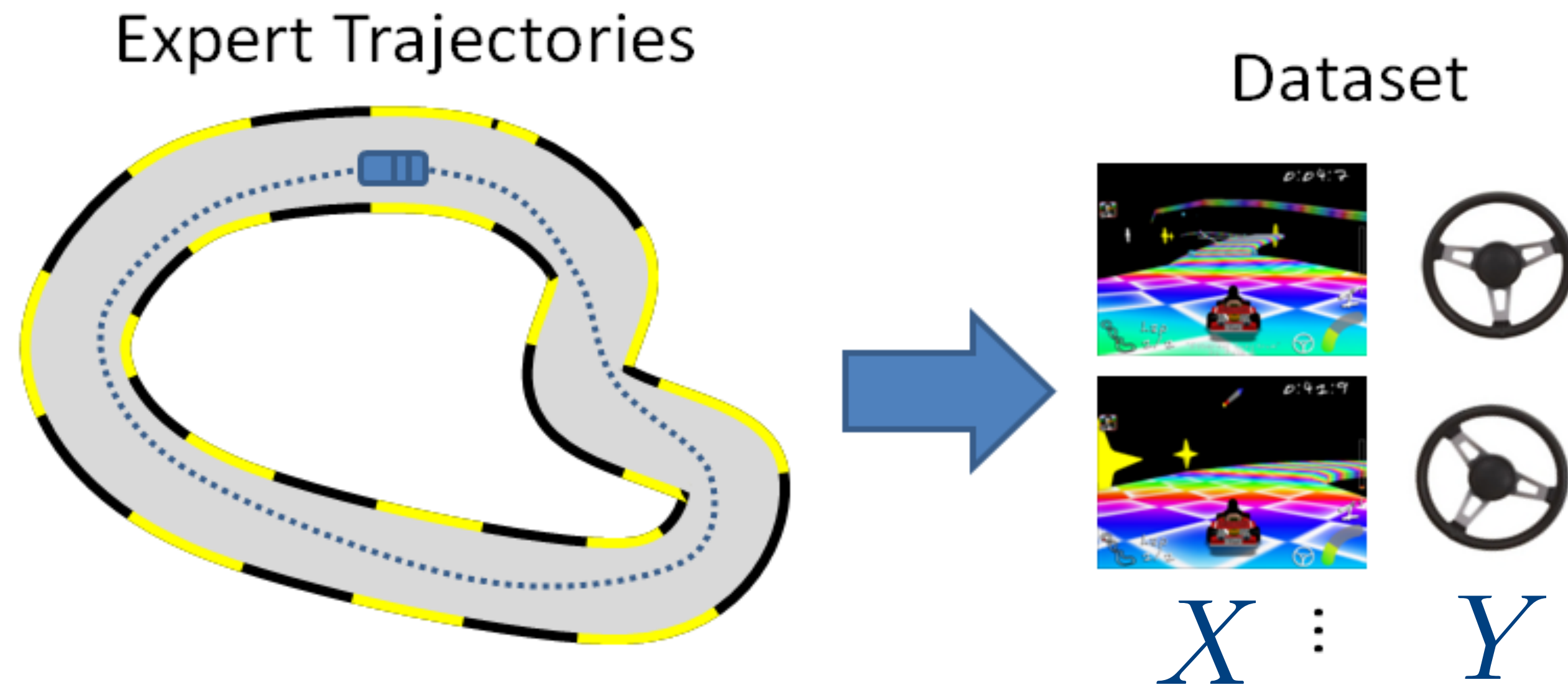


Steering Angle
in $[-1, 1]$

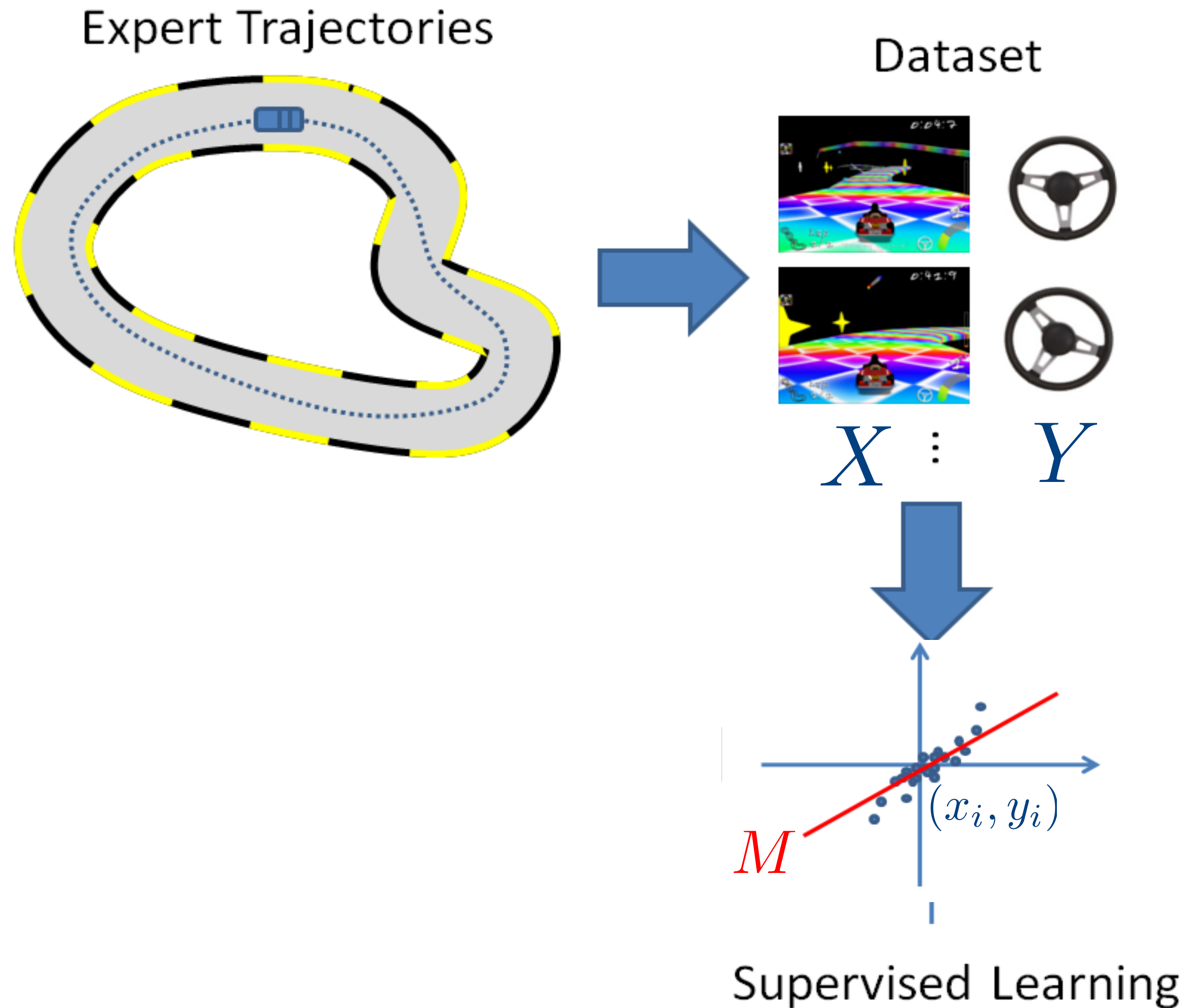
Supervised Learning Approach: Behavior Cloning



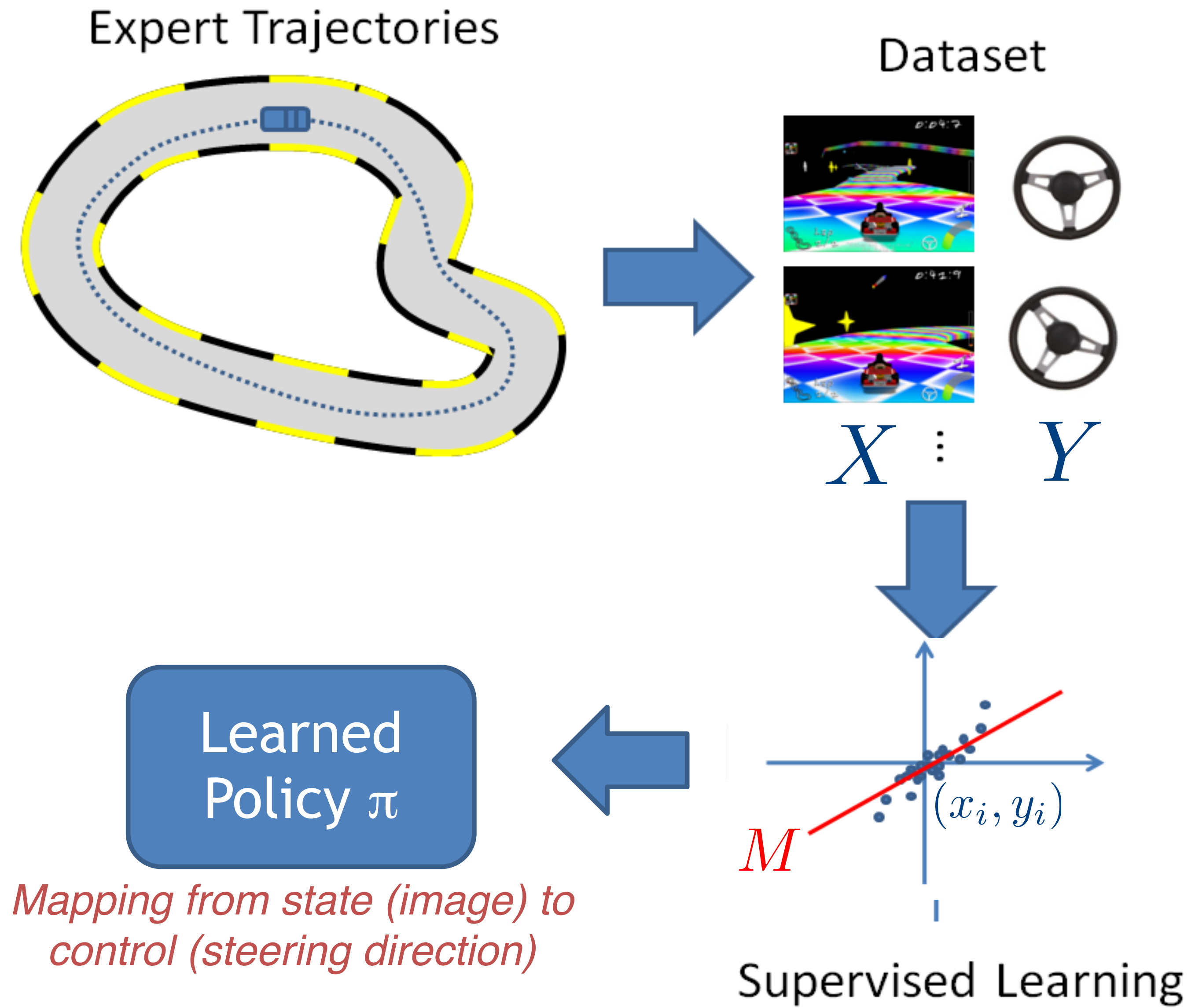
Supervised Learning Approach: Behavior Cloning



Supervised Learning Approach: Behavior Cloning



Supervised Learning Approach: Behavior Cloning









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

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

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

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

Ground truth reward $r(s, a) \in [0, 1]$ is unknown; assume expert is a near optimal policy π^\star

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

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

Ground truth reward $r(s, a) \in [0, 1]$ is unknown; assume expert is a near optimal policy π^\star

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$

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

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

Ground truth reward $r(s, a) \in [0, 1]$ is unknown; assume expert is a near optimal policy π^\star

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$

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

Let's formalize the Behavior Cloning algorithm

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

BC with Maximum Likelihood Estimation (MLE):

Let's formalize the Behavior Cloning algorithm

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

BC with Maximum Likelihood Estimation (MLE):

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Let's formalize the Behavior Cloning algorithm

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

BC with Maximum Likelihood Estimation (MLE):

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

(We can reduce it to other supervised learning oracles such as classification, regression)

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Assumption: Π is discrete, and realizable, i.e., $\pi^* \in \Pi$

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Assumption: Π is discrete, and realizable, i.e., $\pi^* \in \Pi$

Step 1: Supervised learning (MLE) guarantee (see the book for reference to the classic MLE analysis):

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Assumption: Π is discrete, and realizable, i.e., $\pi^* \in \Pi$

Step 1: Supervised learning (MLE) guarantee (see the book for reference to the classic MLE analysis):

Theorem [MLE Guarantee] With probability at least $1 - \delta$, we have:

$$\mathbb{E}_{s \sim d^{\pi^*}} \left\| \hat{\pi}(\cdot | s) - \pi^*(\cdot | s) \right\|_{tv} \leq \sqrt{\frac{\ln(|\Pi|/\delta)}{M}}$$

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Assumption: Π is discrete, and realizable, i.e., $\pi^* \in \Pi$

Step 1: Supervised learning (MLE) guarantee (see the book for reference to the classic MLE analysis):

Theorem [MLE Guarantee] With probability at least $1 - \delta$, we have:

$$\mathbb{E}_{s \sim d^{\pi^*}} \left\| \hat{\pi}(\cdot | s) - \pi^*(\cdot | s) \right\|_{tv} \leq \sqrt{\frac{\ln(|\Pi|/\delta)}{M}}$$

This $1/\sqrt{M}$ rate should be expected:
no training and testing mismatch at this stage!

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Assumption: Π is discrete, and realizable, i.e., $\pi^* \in \Pi$

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Assumption: Π is discrete, and realizable, i.e., $\pi^* \in \Pi$

Step 2: Transfer supervised learning error to policy's performance gap

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Assumption: Π is discrete, and realizable, i.e., $\pi^* \in \Pi$

Step 2: Transfer supervised learning error to policy's performance gap

Theorem [BC Sample Complexity] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

$$V^{\pi^*} - V^{\hat{\pi}} \leq \frac{2}{(1 - \gamma)^2} \underbrace{\sqrt{\frac{\ln(|\Pi|/\delta)}{M}}}_{\text{MLE error}}$$

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Assumption: Π is discrete, and realizable, i.e., $\pi^* \in \Pi$

Step 2: Transfer supervised learning error to policy's performance gap

Theorem [BC Sample Complexity] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

$$V^{\pi^*} - V^{\hat{\pi}} \leq \frac{2}{(1 - \gamma)^2} \underbrace{\sqrt{\frac{\ln(|\Pi|/\delta)}{M}}}_{\text{MLE error}}$$

Note that $1/(1 - \gamma)^2$ quadratic dependency on effective horizon

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Theorem [BC Sample Complexity] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

$$V^{\pi^*} - V^{\hat{\pi}} \leq \frac{2}{(1 - \gamma)^2} \underbrace{\sqrt{\frac{\ln(|\Pi|/\delta)}{M}}}_{\text{MLE error}}$$

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Theorem [BC Sample Complexity] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

$$V^{\pi^*} - V^{\hat{\pi}} \leq \frac{2}{(1 - \gamma)^2} \underbrace{\sqrt{\frac{\ln(|\Pi|/\delta)}{M}}}_{\text{MLE error}}$$

$$(1 - \gamma)(V^* - V^{\hat{\pi}}) = \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot|s)} A^{\hat{\pi}}(s, a)$$

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Theorem [BC Sample Complexity] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

$$V^{\pi^*} - V^{\hat{\pi}} \leq \frac{2}{(1 - \gamma)^2} \underbrace{\sqrt{\frac{\ln(|\Pi|/\delta)}{M}}}_{\text{MLE error}}$$

$$\begin{aligned} (1 - \gamma)(V^* - V^{\hat{\pi}}) &= \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot|s)} A^{\hat{\pi}}(s, a) \\ &= \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot|s)} A^{\hat{\pi}}(s, a) - \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \hat{\pi}(\cdot|s)} A^{\hat{\pi}}(s, a) \end{aligned}$$

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Theorem [BC Sample Complexity] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

$$V^{\pi^*} - V^{\hat{\pi}} \leq \underbrace{\frac{2}{(1-\gamma)^2} \sqrt{\frac{\ln(|\Pi|/\delta)}{M}}}_{\text{MLE error}}$$

$$\begin{aligned} (1-\gamma)(V^* - V^{\hat{\pi}}) &= \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot|s)} A^{\hat{\pi}}(s, a) \\ &= \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot|s)} A^{\hat{\pi}}(s, a) - \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \hat{\pi}(\cdot|s)} A^{\hat{\pi}}(s, a) \\ &\leq \mathbb{E}_{s \sim d^{\pi^*}} \frac{1}{1-\gamma} \left\| \pi^*(\cdot|s) - \hat{\pi}(\cdot|s) \right\|_1 \end{aligned}$$

Analysis

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=1}^M \ln \pi(a_i^* | s_i^*)$$

Theorem [BC Sample Complexity] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

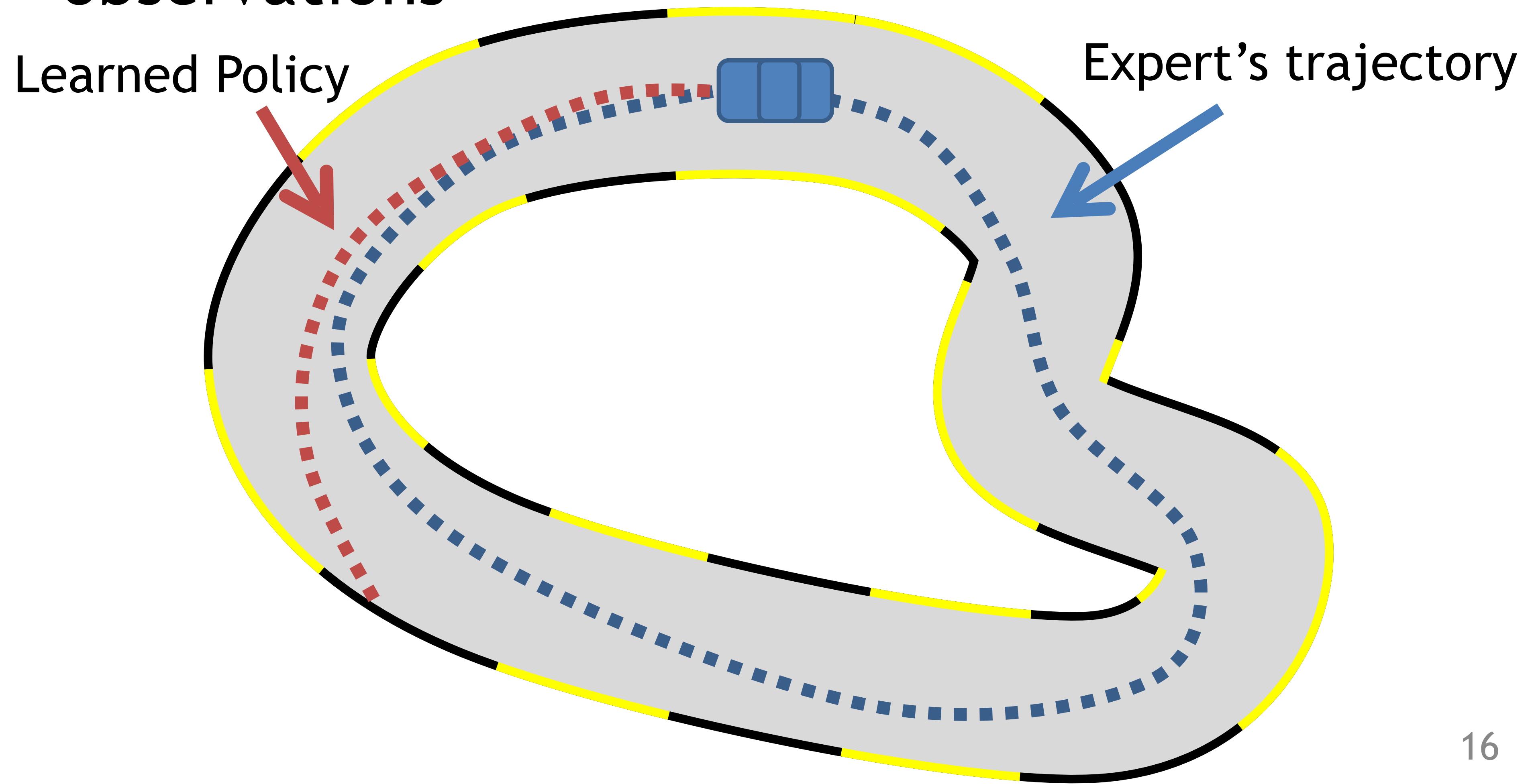
$$V^{\pi^*} - V^{\hat{\pi}} \leq \underbrace{\frac{2}{(1-\gamma)^2} \sqrt{\frac{\ln(|\Pi|/\delta)}{M}}}_{\text{MLE error}}$$

$$\begin{aligned} (1-\gamma)(V^* - V^{\hat{\pi}}) &= \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot|s)} A^{\hat{\pi}}(s, a) \\ &= \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*(\cdot|s)} A^{\hat{\pi}}(s, a) - \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \hat{\pi}(\cdot|s)} A^{\hat{\pi}}(s, a) \\ &\leq \mathbb{E}_{s \sim d^{\pi^*}} \frac{1}{1-\gamma} \left\| \pi^*(\cdot|s) - \hat{\pi}(\cdot|s) \right\|_1 \\ &\leq \frac{2}{1-\gamma} \mathbb{E}_{s \sim d^{\pi^*}} \left\| \pi^*(\cdot|s) - \hat{\pi}(\cdot|s) \right\|_{tv} \end{aligned}$$

What could go wrong?

[Pomerleau89, Daume09]

- Predictions affect future inputs/ observations



Distribution Shift: Intuitive Explanation

Let's just focus on finite horizon (H) and deterministic policies here:

$$\mathbb{E}_{s \sim d_h^{\pi^*}} \hat{\pi}(s) \neq \pi^*(s) \leq \epsilon, \forall h$$

An Autonomous Land Vehicle In A Neural Network *[Pomerleau, NIPS '88]*



An Autonomous Land Vehicle In A Neural Network [Pomerleau, NIPS '88]



“If the network is not presented with sufficient variability in its training exemplars to cover the conditions it is likely to encounter...[it] will perform poorly”

A potential Fix



A potential Fix



Let's roll out our policy in the real world, and compare our trajectories to the expert's trajectories, and then refine our learned model.

The Hybrid Imitation Learning Setting:

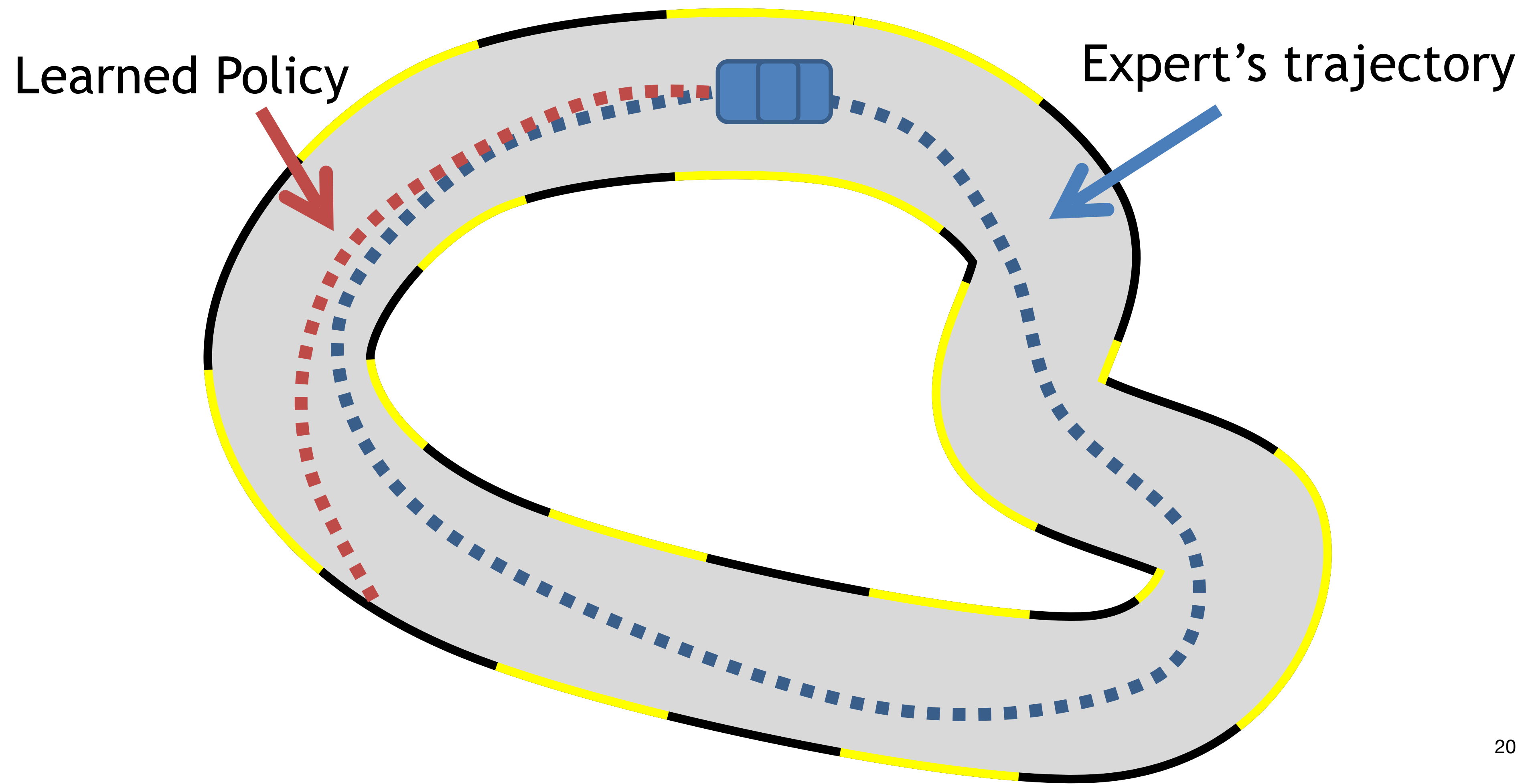
Recall BC: we only use offline expert data—no interaction with the environment

Hybrid setting: offline expert data + simulator (e.g., known transition P)

The Hybrid Imitation Learning Setting:

Recall BC: we only use offline expert data—no interaction with the environment

Hybrid setting: offline expert data + simulator (e.g., known transition P)



Let's formalize the Hybrid Setting

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

Ground truth reward $r(s, a) \in [0, 1]$ is unknown; assume expert is the optimal policy π^\star

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$

Let's formalize the Hybrid Setting

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

Ground truth reward $r(s, a) \in [0, 1]$ is unknown; assume expert is the optimal policy π^\star

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$

This time, we have a known transition P (but we cannot plan because r is unknown)

Let's formalize the Hybrid Setting

Discounted infinite horizon MDP $\mathcal{M} = \{S, A, \gamma, r, P, \rho, \pi^\star\}$

Ground truth reward $r(s, a) \in [0, 1]$ is unknown; assume expert is the optimal policy π^\star

We have a dataset $\mathcal{D} = (s_i^\star, a_i^\star)_{i=1}^M \sim d^{\pi^\star}$

This time, we have a known transition P (but we cannot plan because r is unknown)

Key Q: can we do better than offline IL Behavior Cloning (statistically at least—assuming infinite computation power)?

Key Idea: distribution matching

Integral probability metric (IPM)

Key Idea: distribution matching

Integral probability metric (IPM)

Metric measures the divergence between two distributions

Key Idea: distribution matching

Integral probability metric (IPM)

Metric measures the divergence between two distributions

Given a discriminator class: $\mathcal{F} = \{f : X \mapsto \mathbb{R}\}$, and two distributions p_1 and p_2

Key Idea: distribution matching

Integral probability metric (IPM)

Metric measures the divergence between two distributions

Given a discriminator class: $\mathcal{F} = \{f : X \mapsto \mathbb{R}\}$, and two distributions p_1 and p_2

$$\text{IPM}_{\mathcal{F}}(p_1, p_2) = \max_{f \in \mathcal{F}} \left[\mathbb{E}_{x \sim p_1} f(x) - \mathbb{E}_{x \sim p_2} f(x) \right]$$

Key Idea: distribution matching

Integral probability metric (IPM)

Metric measures the divergence between two distributions

Given a discriminator class: $\mathcal{F} = \{f : X \mapsto \mathbb{R}\}$, and two distributions p_1 and p_2

$$\text{IPM}_{\mathcal{F}}(p_1, p_2) = \max_{f \in \mathcal{F}} \left[\mathbb{E}_{x \sim p_1} f(x) - \mathbb{E}_{x \sim p_2} f(x) \right]$$

$$\mathcal{F} = \{f : \|f\|_{\infty} \leq 1\} \Rightarrow \text{IPM}_{\mathcal{F}}(p_1, p_2) := \|p_1 - p_2\|_{tv}$$

Key Idea: distribution matching

Integral probability metric (IPM)

Metric measures the divergence between two distributions

Given a discriminator class: $\mathcal{F} = \{f : X \mapsto \mathbb{R}\}$, and two distributions p_1 and p_2

$$\text{IPM}_{\mathcal{F}}(p_1, p_2) = \max_{f \in \mathcal{F}} \left[\mathbb{E}_{x \sim p_1} f(x) - \mathbb{E}_{x \sim p_2} f(x) \right]$$

$$\mathcal{F} = \{f : \|f\|_{\infty} \leq 1\} \Rightarrow \text{IPM}_{\mathcal{F}}(p_1, p_2) := \|p_1 - p_2\|_{tv}$$

$$\mathcal{F} = \{f : f \text{ is 1-Lipschitz}\} \Rightarrow \text{IPM}_{\mathcal{F}}(p_1, p_2) := \text{wasserstein dis}(p_1, p_2)$$

Algorithm: Distribution Matching TV distance

Consider the Discriminator class: $\mathcal{F} = \{f : \|f\|_\infty \leq 1\}$ (IPM corresponds to TV distance)

Algorithm: Distribution Matching TV distance

Consider the Discriminator class: $\mathcal{F} = \{f : \|f\|_\infty \leq 1\}$ (IPM corresponds to TV distance)

Same assumption as we had in BC:
 Π is discrete, and realizable, i.e., $\pi^\star \in \Pi$

Algorithm: Distribution Matching TV distance

Consider the Discriminator class: $\mathcal{F} = \{f : \|f\|_\infty \leq 1\}$ (IPM corresponds to TV distance)

Same assumption as we had in BC:

Π is discrete, and realizable, i.e., $\pi^\star \in \Pi$

Step 1: for each $\pi \in \Pi$, compute $d^\pi \in \Delta(S \times A)$ (recall P is known);
(This step is computationally inefficient)

Algorithm: Distribution Matching TV distance

Consider the Discriminator class: $\mathcal{F} = \{f : \|f\|_\infty \leq 1\}$ (IPM corresponds to TV distance)

Same assumption as we had in BC:

Π is discrete, and realizable, i.e., $\pi^\star \in \Pi$

Step 1: for each $\pi \in \Pi$, compute $d^\pi \in \Delta(S \times A)$ (recall P is known);
(This step is computationally inefficient)

Step 2: select useful discriminators: for all pair π & π' , with $\pi \neq \pi'$

$$f_{\pi, \pi'} = \arg \max_{f \in \mathcal{F}} \left[\mathbb{E}_{s, a \sim d^\pi} f(s, a) - \mathbb{E}_{s, a \sim d^{\pi'}} f(s, a) \right]$$

Algorithm: Distribution Matching TV distance

Consider the Discriminator class: $\mathcal{F} = \{f : \|f\|_\infty \leq 1\}$ (IPM corresponds to TV distance)

Same assumption as we had in BC:

Π is discrete, and realizable, i.e., $\pi^\star \in \Pi$

Step 1: for each $\pi \in \Pi$, compute $d^\pi \in \Delta(S \times A)$ (recall P is known);
(This step is computationally inefficient)

Step 2: select useful discriminators: for all pair π & π' , with $\pi \neq \pi'$

$$f_{\pi,\pi'} = \arg \max_{f \in \mathcal{F}} \left[\mathbb{E}_{s,a \sim d^\pi} f(s, a) - \mathbb{E}_{s,a \sim d^{\pi'}} f(s, a) \right]$$

Set refined discriminator class $\widetilde{\mathcal{F}} := \{f_{\pi,\pi'} : \pi \text{ \& \ } \pi' \in \Pi, \pi \neq \pi'\}$

Algorithm: Distribution Matching TV distance

Consider the Discriminator class: $\mathcal{F} = \{f : \|f\|_\infty \leq 1\}$ (IPM corresponds to TV distance)

Same assumption as we had in BC:

Π is discrete, and realizable, i.e., $\pi^\star \in \Pi$

Step 1: for each $\pi \in \Pi$, compute $d^\pi \in \Delta(S \times A)$ (recall P is known);
(This step is computationally inefficient)

Step 2: select useful discriminators: for all pair π & π' , with $\pi \neq \pi'$

$$f_{\pi,\pi'} = \arg \max_{f \in \mathcal{F}} \left[\mathbb{E}_{s,a \sim d^\pi} f(s, a) - \mathbb{E}_{s,a \sim d^{\pi'}} f(s, a) \right]$$

Q: what is the size of $\widetilde{\mathcal{F}}$?

Set refined discriminator class $\widetilde{\mathcal{F}} := \{f_{\pi,\pi'} : \pi \text{ \& \ } \pi' \in \Pi, \pi \neq \pi'\}$

Algorithm: Distribution Matching TV distance

Consider the Discriminator class: $\mathcal{F} = \{f : \|f\|_\infty \leq 1\}$ (IPM corresponds to TV distance)

Same assumption as we had in BC:
 Π is discrete, and realizable, i.e., $\pi^\star \in \Pi$

Step 1: for each $\pi \in \Pi$, compute $d^\pi \in \Delta(S \times A)$ (recall P is known);
(This step is computationally inefficient)

Step 2: select useful discriminators: for all pair π & π' , with $\pi \neq \pi'$

$$f_{\pi,\pi'} = \arg \max_{f \in \mathcal{F}} \left[\mathbb{E}_{s,a \sim d^\pi} f(s, a) - \mathbb{E}_{s,a \sim d^{\pi'}} f(s, a) \right]$$

Q: what is the size of $\widetilde{\mathcal{F}}$?

Set refined discriminator class $\widetilde{\mathcal{F}} := \{f_{\pi,\pi'} : \pi \text{ \& \ } \pi' \in \Pi, \pi \neq \pi'\}$

$$\forall \pi \text{ \& \ } \pi' \in \Pi, \|d^\pi - d^{\pi'}\|_{tv} = \max_{f \in \widetilde{\mathcal{F}}} \mathbb{E}_{s,a \sim d^\pi} f(s, a) - \mathbb{E}_{s,a \sim d^{\pi'}} f(s, a)$$

Algorithm: Distribution Matching TV distance

Step 1: for each $\pi \in \Pi$, compute $d^\pi \in \Delta(S \times A)$ (recall P is known);
(This step is computationally inefficient)

Step2: select useful discriminators: for all pair π & π' , with $\pi \neq \pi'$

$$f_{\pi, \pi'} = \arg \max_{f \in \mathcal{F}} \left[\mathbb{E}_{s, a \sim d^\pi} f(s, a) - \mathbb{E}_{s, a \sim d^{\pi'}} f(s, a) \right]$$

Set refined discriminator class $\widetilde{\mathcal{F}} := \{f_{\pi, \pi'} : \pi \text{ \& \ } \pi' \in \Pi, \pi \neq \pi'\}$

Algorithm: Distribution Matching TV distance

Step 1: for each $\pi \in \Pi$, compute $d^\pi \in \Delta(S \times A)$ (recall P is known);
(This step is computationally inefficient)

Step2: select useful discriminators: for all pair π & π' , with $\pi \neq \pi'$

$$f_{\pi, \pi'} = \arg \max_{f \in \mathcal{F}} \left[\mathbb{E}_{s, a \sim d^\pi} f(s, a) - \mathbb{E}_{s, a \sim d^{\pi'}} f(s, a) \right]$$

Set refined discriminator class $\widetilde{\mathcal{F}} := \{f_{\pi, \pi'} : \pi \text{ \& \ } \pi' \in \Pi, \pi \neq \pi'\}$

Step 3: Select a policy using expert dataset $\mathcal{D} = \{s_i^\star, a_i^\star\}_{i=1}^M$

Algorithm: Distribution Matching TV distance

Step 1: for each $\pi \in \Pi$, compute $d^\pi \in \Delta(S \times A)$ (recall P is known);
(This step is computationally inefficient)

Step2: select useful discriminators: for all pair π & π' , with $\pi \neq \pi'$

$$f_{\pi, \pi'} = \arg \max_{f \in \mathcal{F}} \left[\mathbb{E}_{s, a \sim d^\pi} f(s, a) - \mathbb{E}_{s, a \sim d^{\pi'}} f(s, a) \right]$$

Set refined discriminator class $\widetilde{\mathcal{F}} := \{f_{\pi, \pi'} : \pi \text{ \& \ } \pi' \in \Pi, \pi \neq \pi'\}$

Step 3: Select a policy using expert dataset $\mathcal{D} = \{s_i^\star, a_i^\star\}_{i=1}^M$

$$\hat{\pi} := \arg \min_{\pi \in \Pi} \left[\max_{f \in \widetilde{\mathcal{F}}} \left[\mathbb{E}_{s, a \sim d^\pi} f(s, a) - \frac{1}{M} \sum_{i=1}^M f(s_i^\star, a_i^\star) \right] \right]$$

Theorem: Distribution Matching TV distance

Theorem [Dis-match w/ TV dist] With probability at least $1 - \delta$, our algorithm finds a policy $\hat{\pi}$, s.t.,

$$V^{\pi^*} - V^{\hat{\pi}} \leq \mathcal{O} \left(\frac{1}{1 - \gamma} \sqrt{\frac{\ln(|\Pi|/\delta)}{M}} \right)$$

Theorem: Distribution Matching TV distance

Theorem [Dis-match w/ TV dist] With probability at least $1 - \delta$, our algorithm finds a policy $\hat{\pi}$, s.t.,

$$V^{\pi^*} - V^{\hat{\pi}} \leq \mathcal{O} \left(\frac{1}{1 - \gamma} \sqrt{\frac{\ln(|\Pi|/\delta)}{M}} \right)$$

1. Key step is to prove: $\| d^{\hat{\pi}} - d^{\pi^*} \|_{tv} \leq \sqrt{\frac{\ln(|\Pi|/\delta)}{M}}$

Theorem: Distribution Matching TV distance

Theorem [Dis-match w/ TV dist] With probability at least $1 - \delta$, our algorithm finds a policy $\hat{\pi}$, s.t.,

$$V^{\pi^*} - V^{\hat{\pi}} \leq \mathcal{O} \left(\frac{1}{1 - \gamma} \sqrt{\frac{\ln(|\Pi|/\delta)}{M}} \right)$$

1. Key step is to prove: $\| d^{\hat{\pi}} - d^{\pi^*} \|_{tv} \leq \sqrt{\frac{\ln(|\Pi|/\delta)}{M}}$

2. For performance: $V^{\pi^*} - V^{\pi} \leq \frac{1}{1 - \gamma} \left[\mathbb{E}_{s,a \sim d^{\pi^*}} r(s, a) - \mathbb{E}_{s,a \sim d^{\hat{\pi}}} r(s, a) \right]$

Theorem: Distribution Matching TV distance

Theorem [Dis-match w/ TV dist] With probability at least $1 - \delta$, our algorithm finds a policy $\hat{\pi}$, s.t.,

$$V^{\pi^*} - V^{\hat{\pi}} \leq \mathcal{O} \left(\frac{1}{1 - \gamma} \sqrt{\frac{\ln(|\Pi|/\delta)}{M}} \right)$$

1. Key step is to prove: $\| d^{\hat{\pi}} - d^{\pi^*} \|_{tv} \leq \sqrt{\frac{\ln(|\Pi|/\delta)}{M}}$

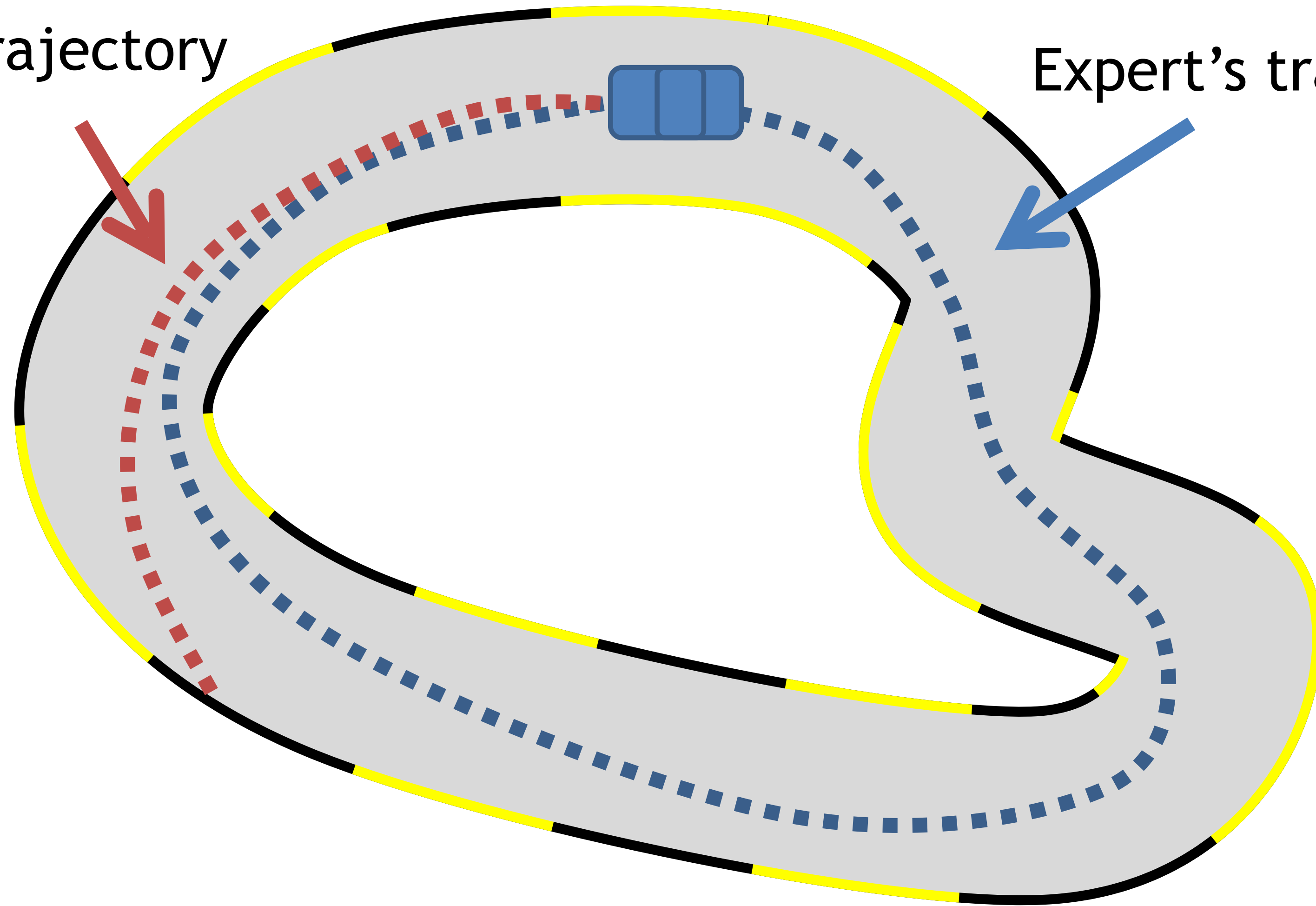
2. For performance: $V^{\pi^*} - V^{\hat{\pi}} \leq \frac{1}{1 - \gamma} \left[\mathbb{E}_{s,a \sim d^{\pi^*}} r(s, a) - \mathbb{E}_{s,a \sim d^{\hat{\pi}}} r(s, a) \right]$

Theorem [Offline BC] With probability at least $1 - \delta$, BC returns a policy $\hat{\pi}$:

$$V^{\pi^*} - V^{\hat{\pi}} = \mathcal{O} \left(\frac{1}{(1 - \gamma)^2} \sqrt{\frac{\ln(|\Pi|/\delta)}{M}} \right)$$

A Policy's trajectory

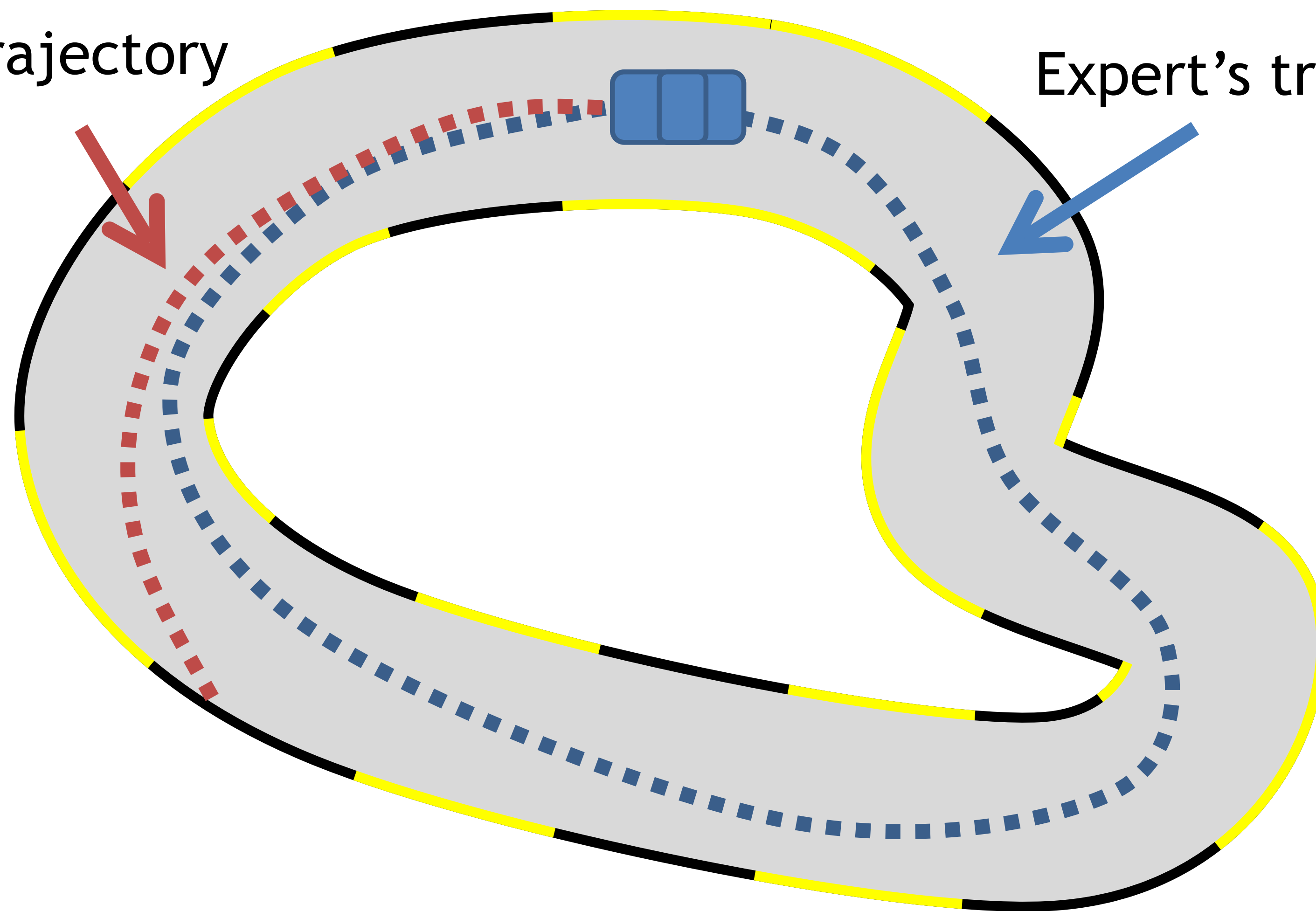
Expert's trajectory



d^π : generator that generates state-action pairs

A Policy's trajectory

Expert's trajectory

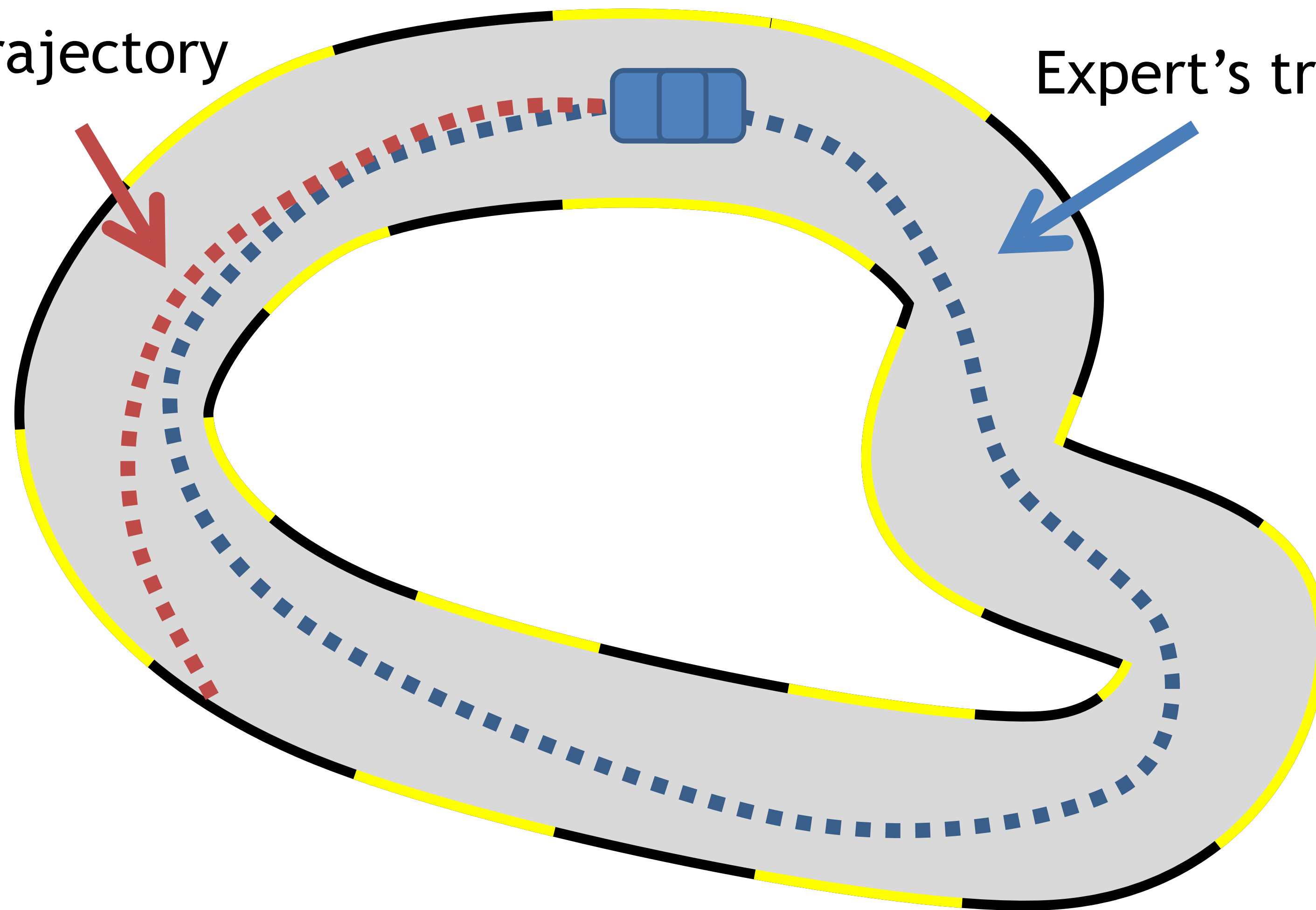


d^π : generator that generates state-action pairs

d^{π^*} : Ground truth state-action distribution (we have samples from it)

A Policy's trajectory

Expert's trajectory



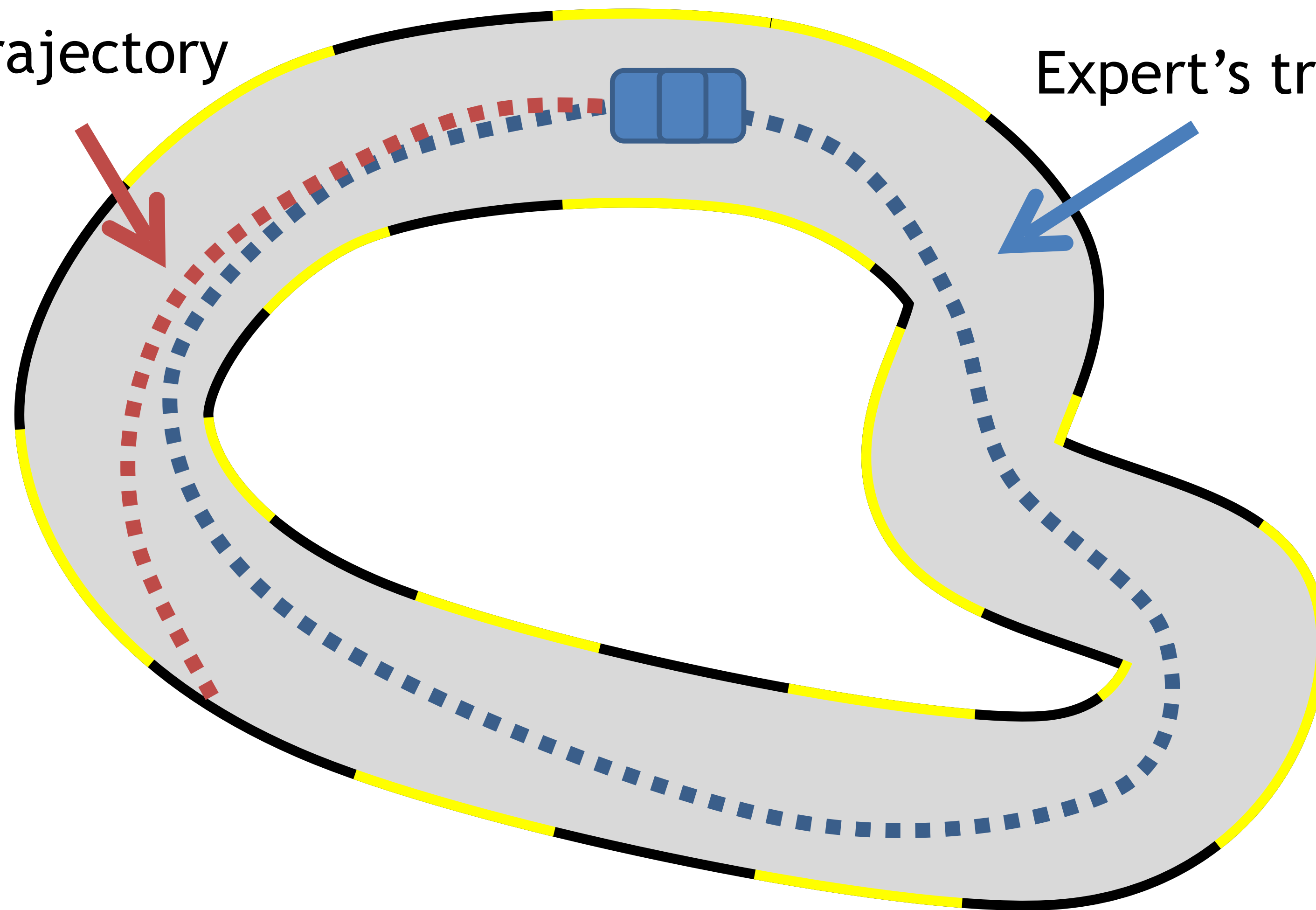
d^π : generator that generates state-action pairs

d^{π^*} : Ground truth state-action distribution (we have samples from it)

\mathcal{F} : discriminators which distinguish red and blue

A Policy's trajectory

Expert's trajectory



Conclusion:

1. Offline RL: only use offline expert data

BC is **simple and easy to implement**, has reasonable guarantees; but the **quadratic dependency on horizon** could cause real problems

Conclusion:

1. Offline RL: only use offline expert data

BC is **simple and easy to implement**, has reasonable guarantees; but the **quadratic dependency on horizon** could cause real problems

2. Hybrid RL: offline expert data + known transition (simulator)

Statistically, Distribution-matching has **linear dependency on horizon**, but the algorithm is **computationally inefficient**

Conclusion:

1. Offline RL: only use offline expert data

BC is **simple and easy to implement**, has reasonable guarantees; but the **quadratic dependency on horizon** could cause real problems

2. Hybrid RL: offline expert data + known transition (simulator)

Statistically, Distribution-matching has **linear dependency on horizon**, but the algorithm is **computationally inefficient**

Take home message:

There is a **provable statistical benefit from the hybrid setting!**

Ps: the distribution matching algorithm is very new (it was discovered when I was writing the book chapter...)