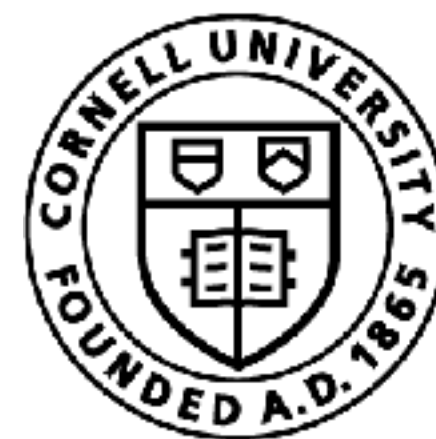


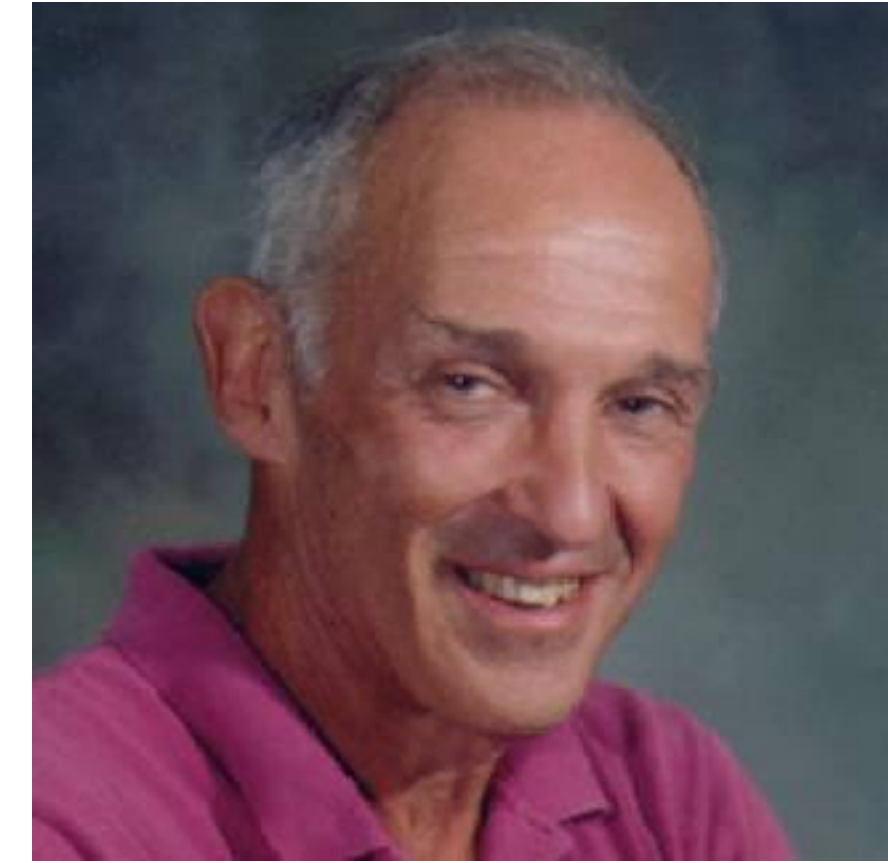
DAgger: Taming Covariate Shift with No Regret

Sanjiban Choudhury



Cornell Bowers CIS
Computer Science

Really cool talk today!



My Six Decade Experience in Visual Imaging

Prof. Don Greenburg

Colloquium at 11:45am, in G01 Gates Hall

RESEARCH NIGHT

Monday
Sept. 16 2024

Panel Discussion
Gates G01 5-6 PM

Poster Session
Gates Atrium 6-7
PM



ACSU

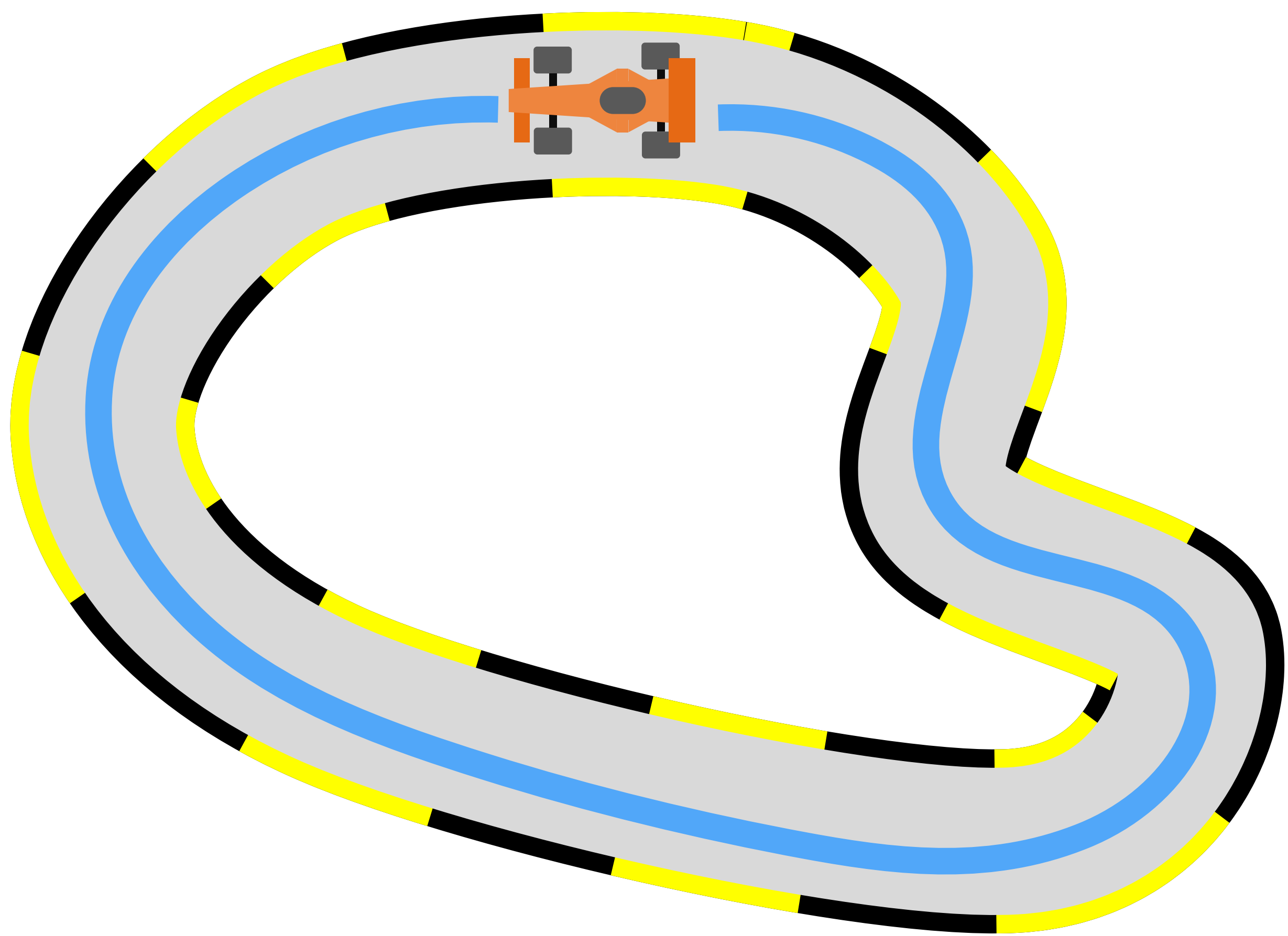
Gates Hall



Today's class

- What is DAGGER?
- DAGGER in the real world
- Why *aggregate* data?

Behavior Cloning

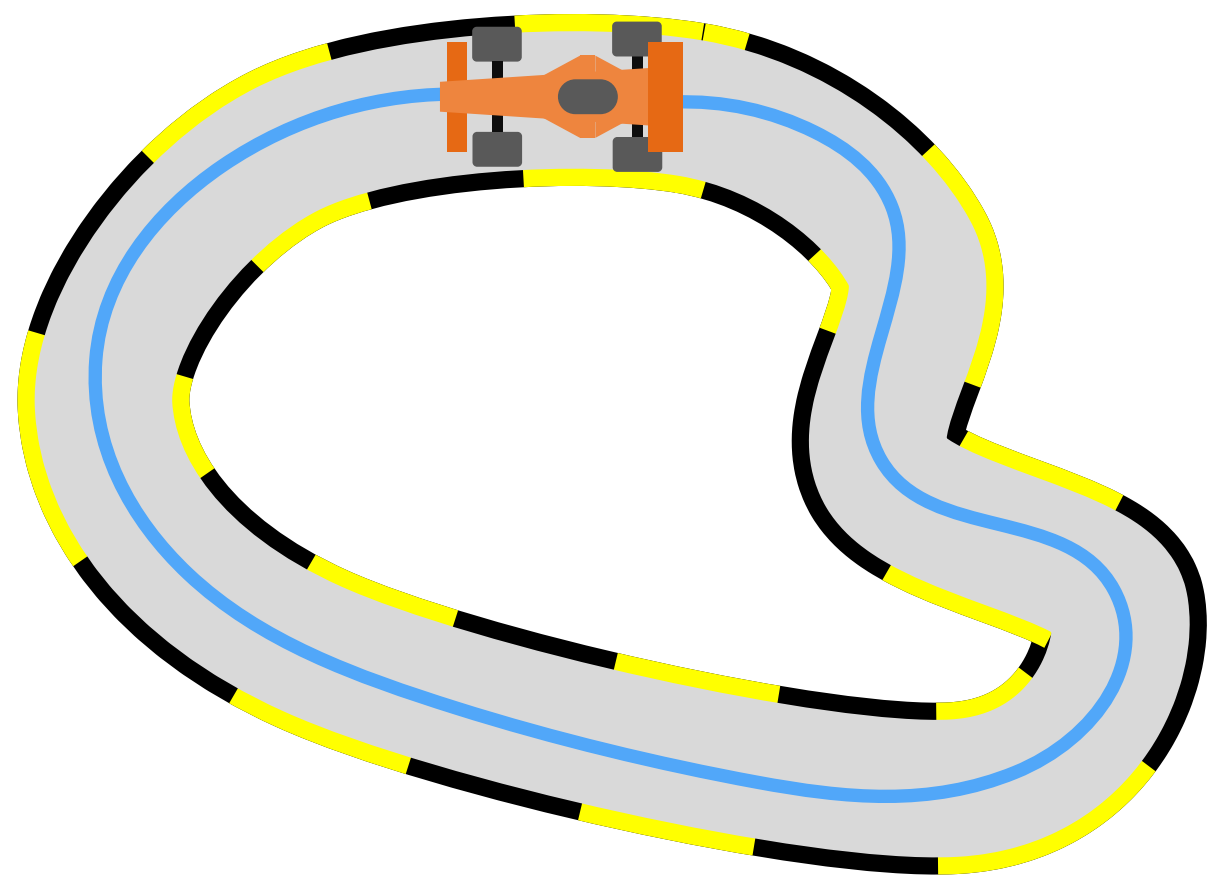


Expert runs away after demonstrations

The Big Problem with BC

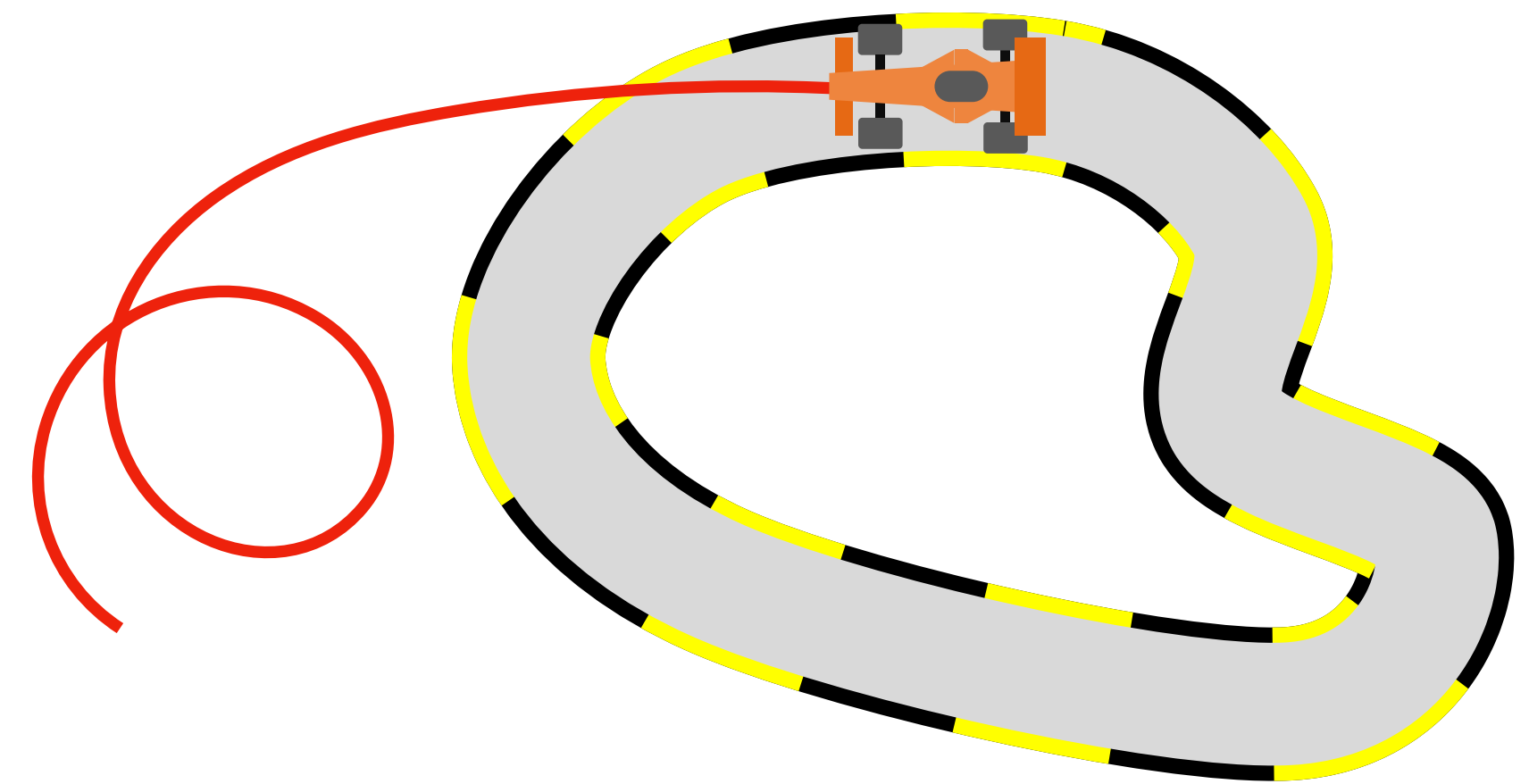
Train

$$\sum_{t=0}^{T-1} \mathbb{E}_{s_t \sim d_t^{\pi^*}} [\ell(s_t, \pi(s_t))]$$



Test

$$\sum_{t=0}^{T-1} \mathbb{E}_{s_t \sim d_t^{\pi}} [\ell(s_t, \pi(s_t))]$$



Goal: Bound **on-policy** loss

$$\sum_{t=0}^{T-1} \mathbb{E}_{s_t \sim d_t^\pi} [\ell(s_t, \pi(s_t))]$$

If we can bound above, we can bound performance difference

$$J(\pi) - J(\pi^*)$$



What if we interactively
queried the expert on
states the learner visits?

DAGGER: A **meta-algorithm** for imitation learning

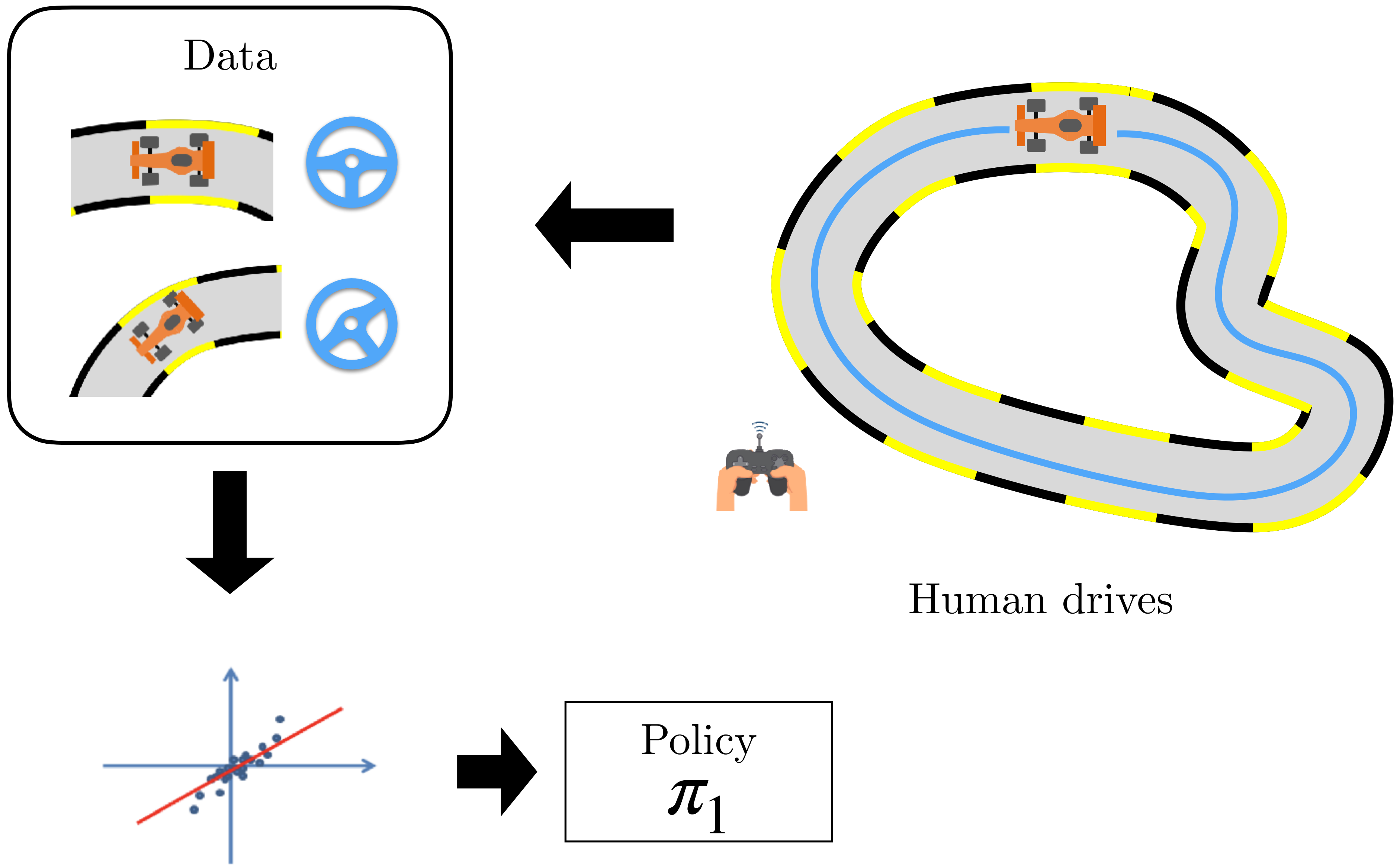
A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning

Stéphane Ross
Robotics Institute
Carnegie Mellon University
Pittsburgh, PA 15213, USA
stephaneross@cmu.edu

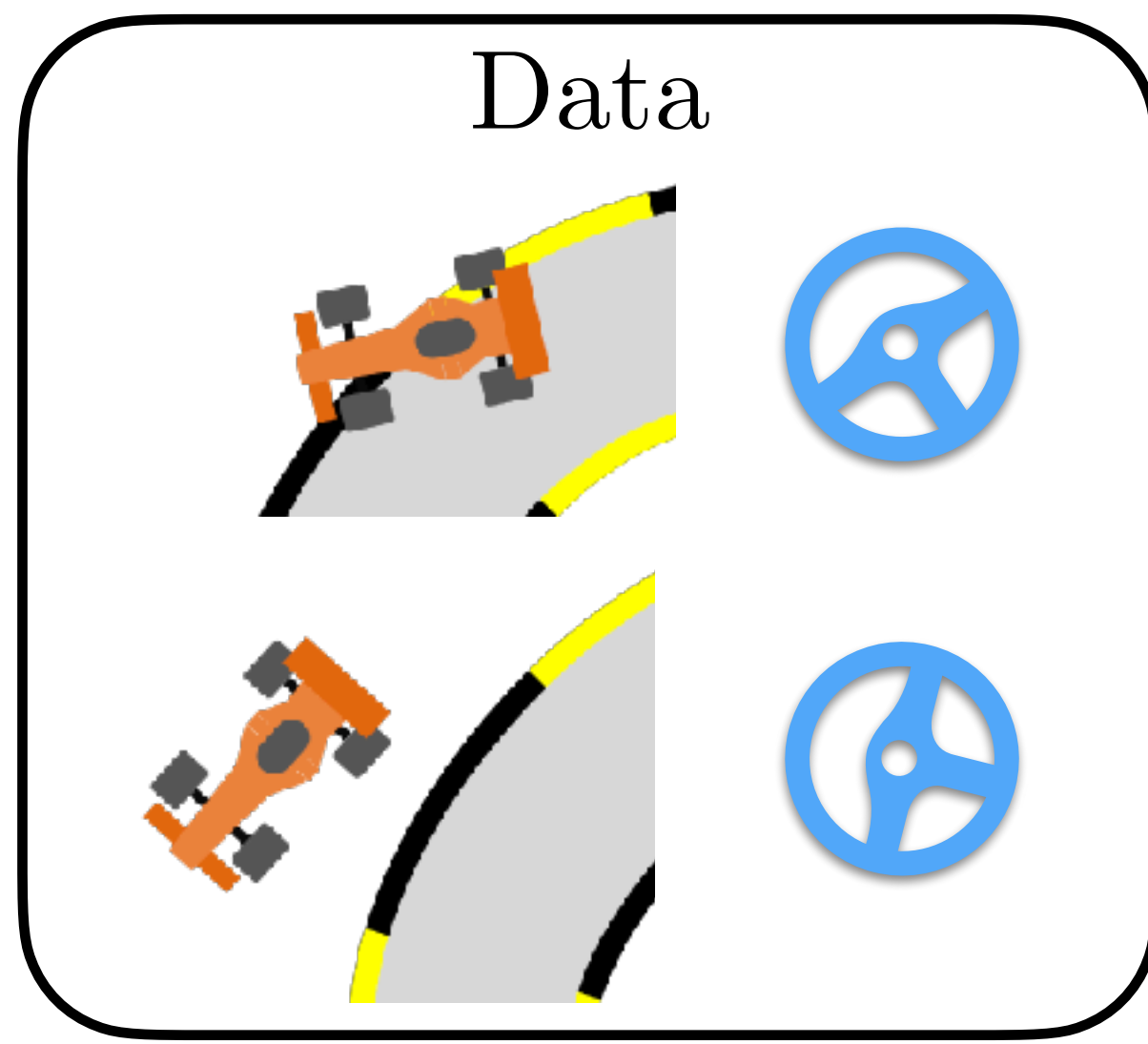
Geoffrey J. Gordon
Machine Learning Department
Carnegie Mellon University
Pittsburgh, PA 15213, USA
ggordon@cs.cmu.edu

J. Andrew Bagnell
Robotics Institute
Carnegie Mellon University
Pittsburgh, PA 15213, USA
dbagnell@ri.cmu.edu

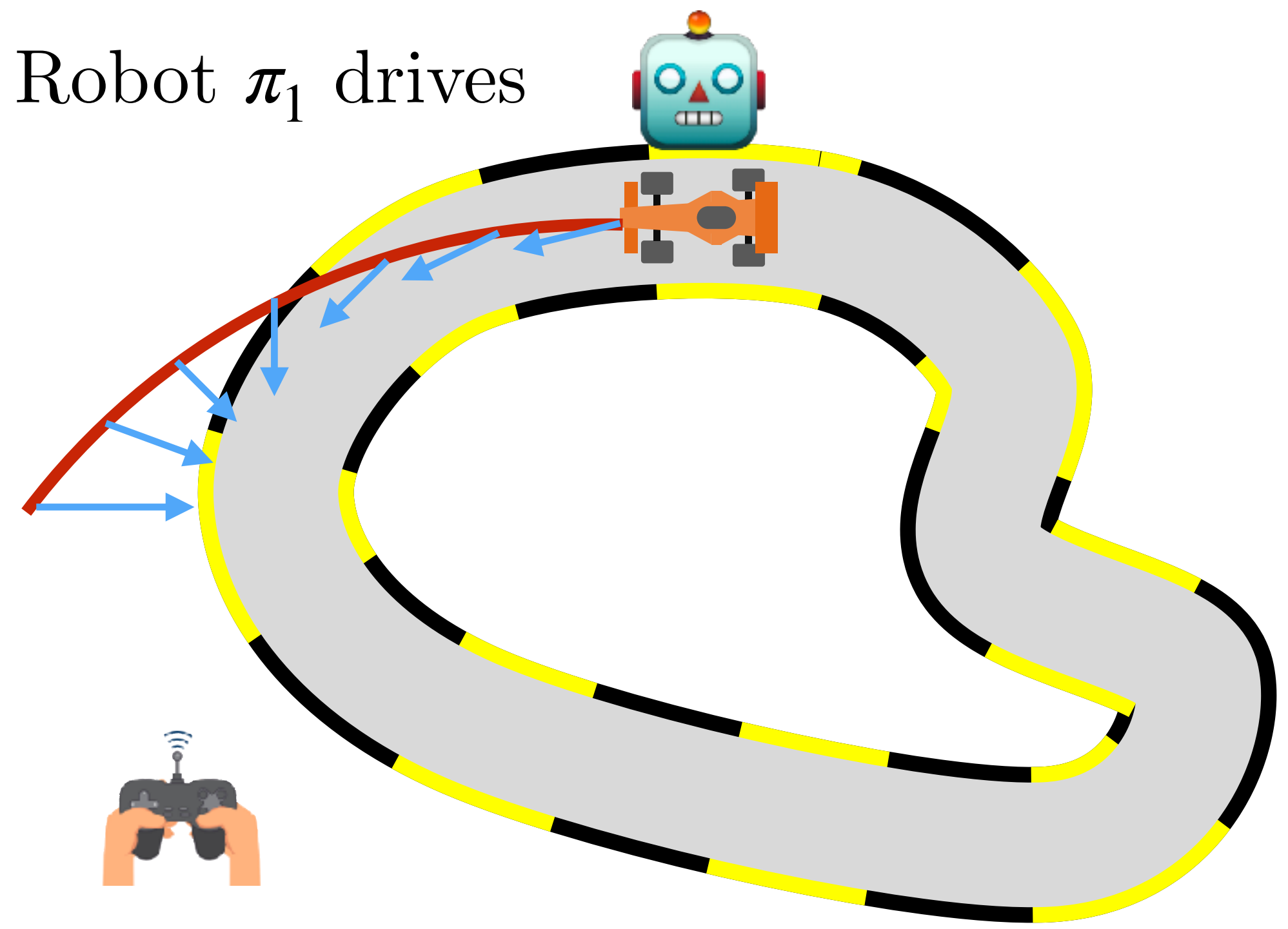
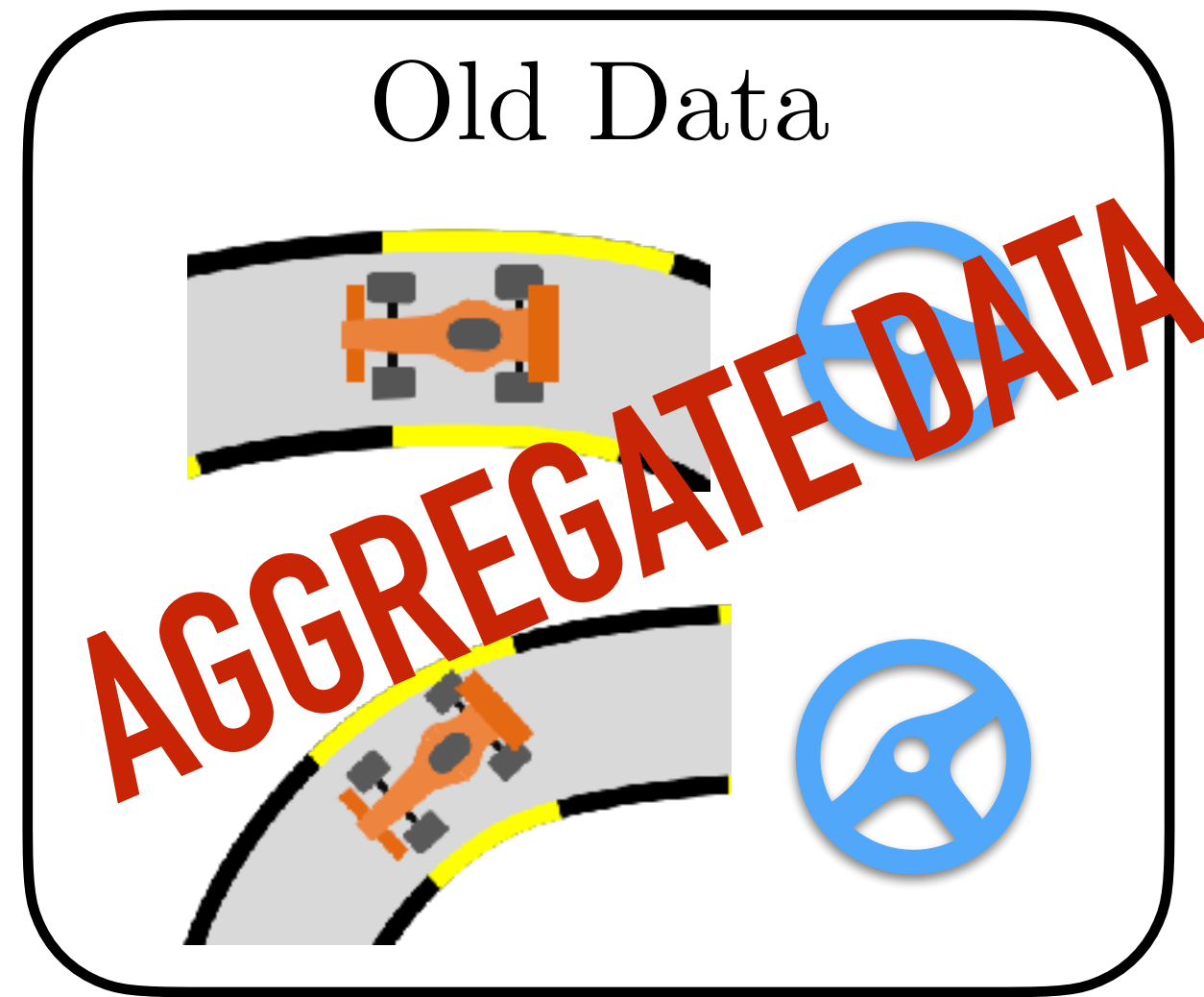
DAgger: Initializations



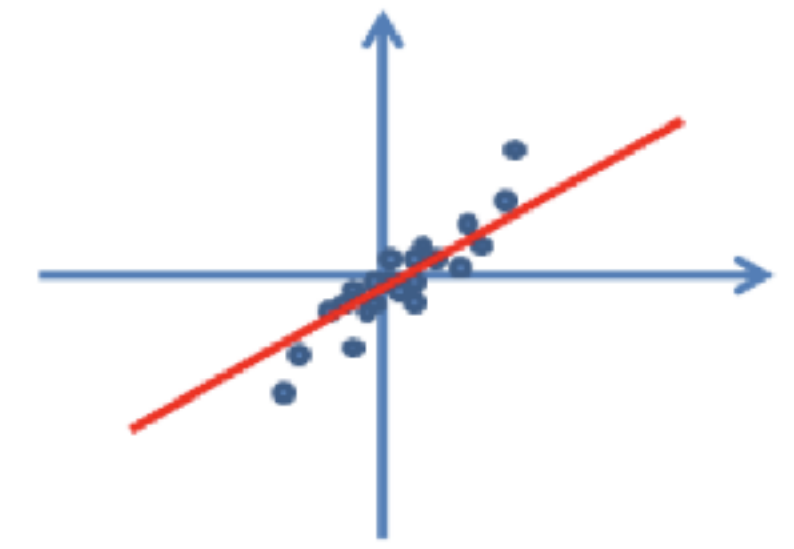
Dagger: Iteration 1



+

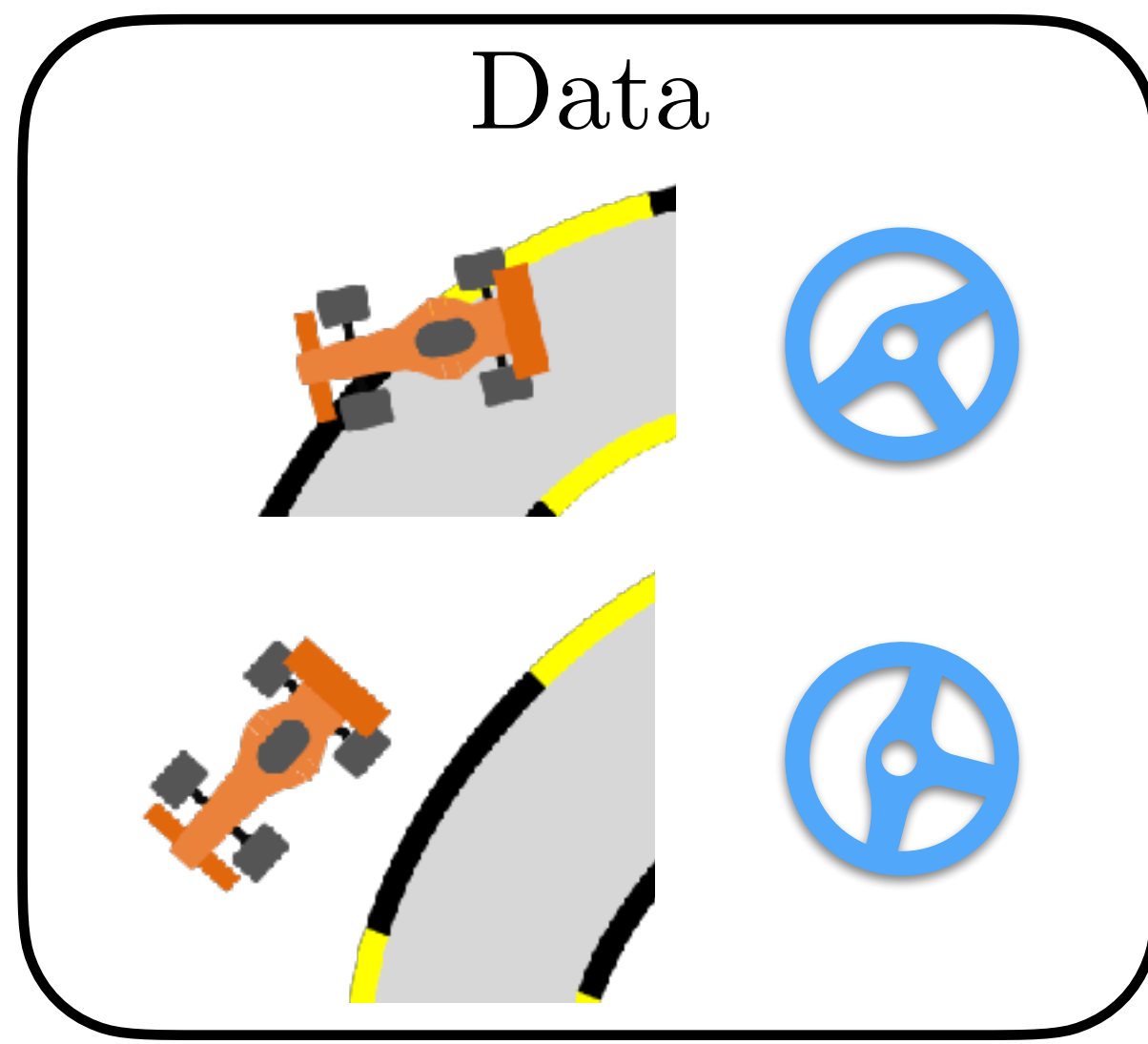


Human corrects!

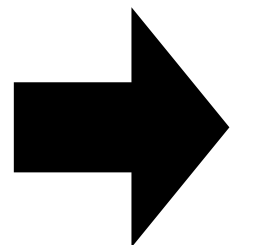
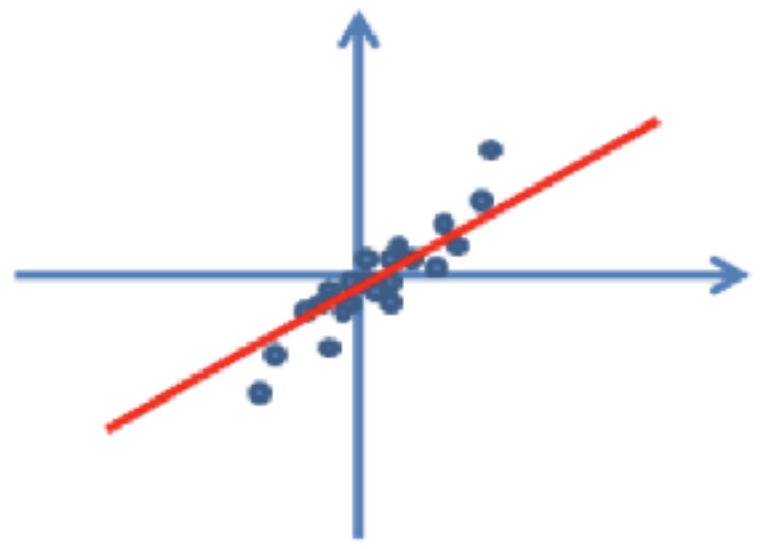
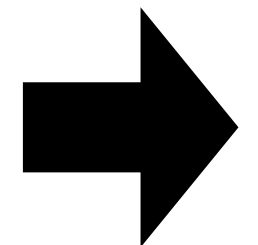
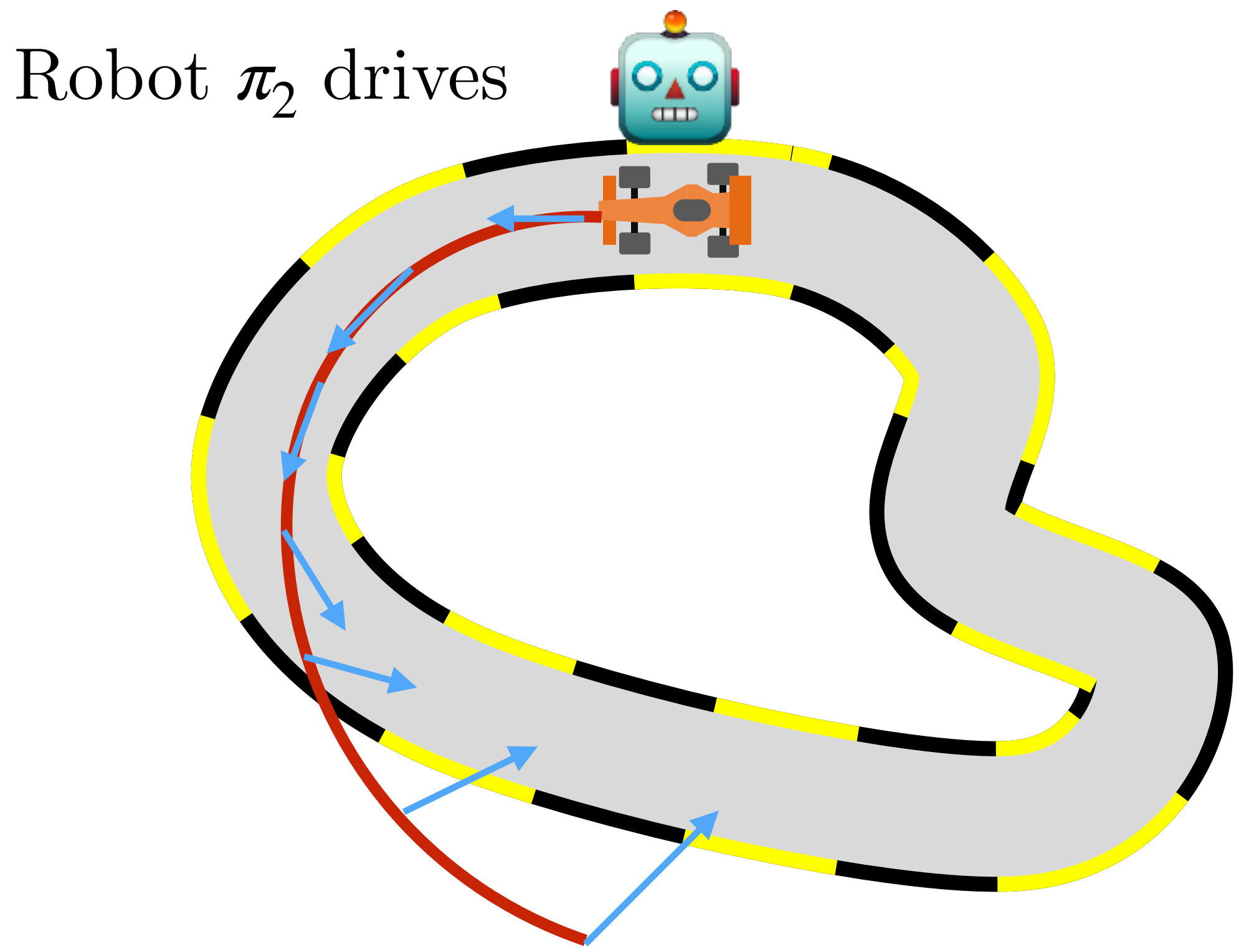
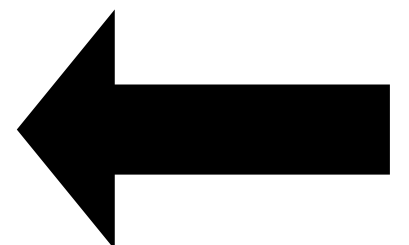
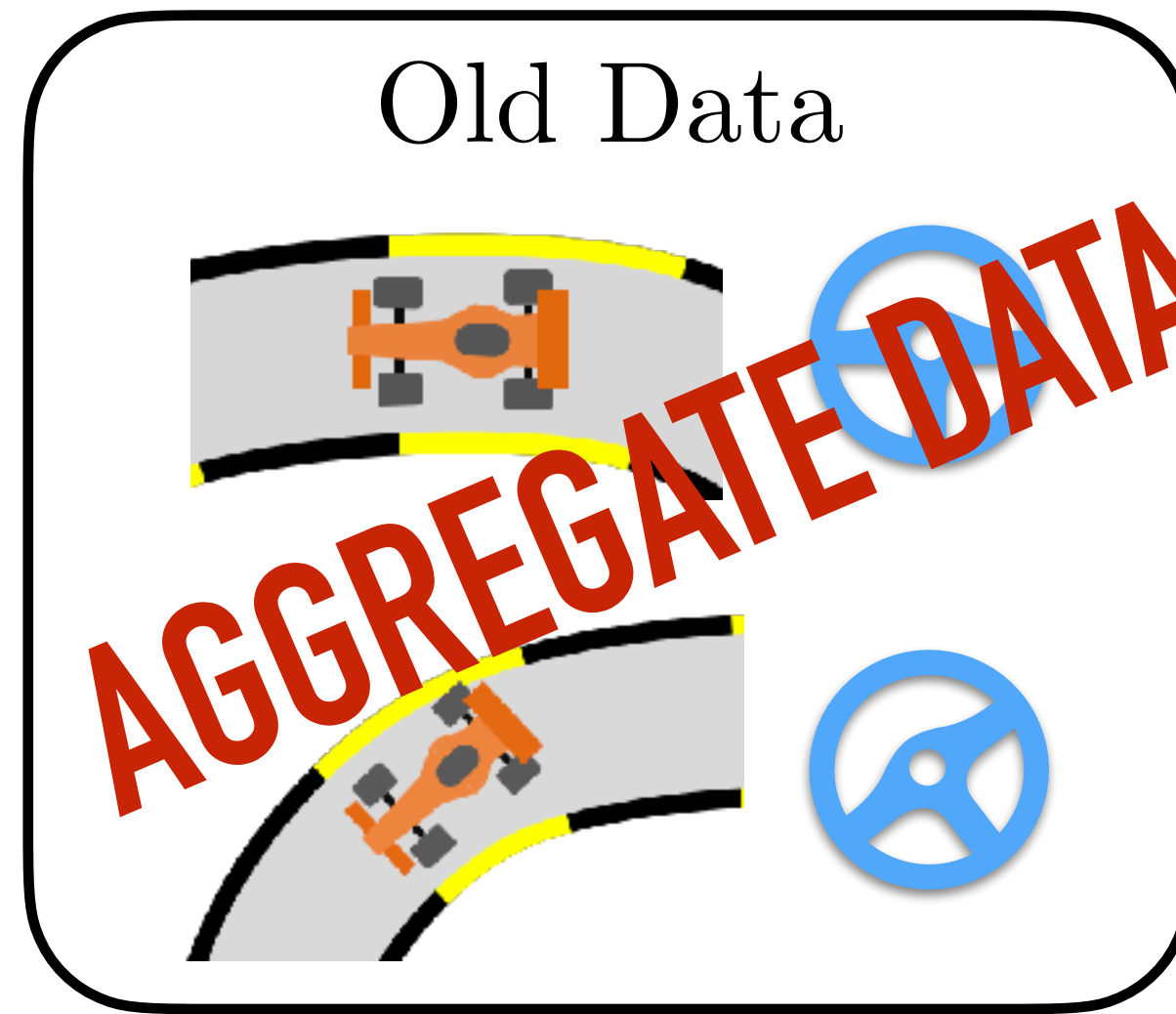


Policy π_2

Dagger: Iteration 2



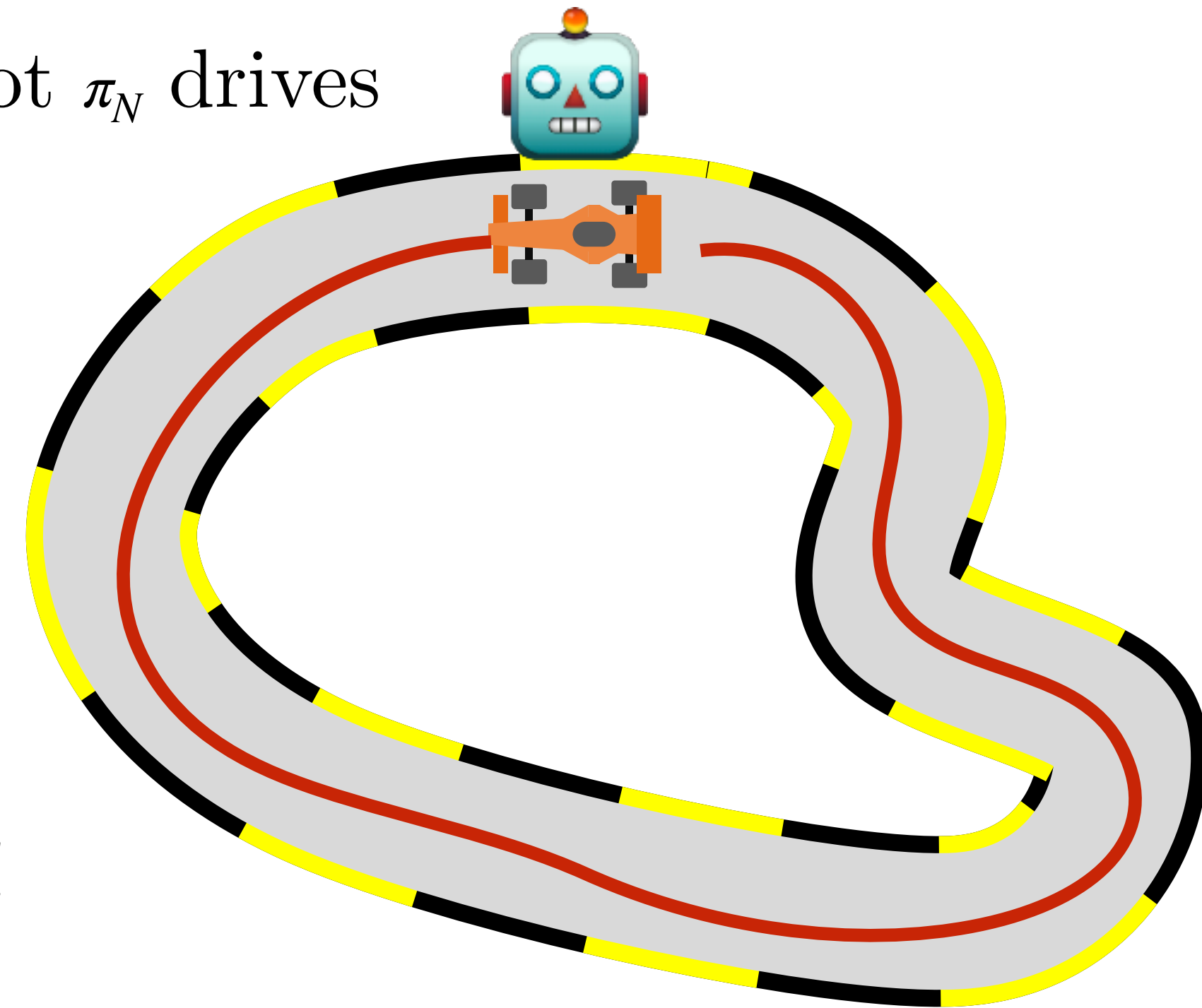
+



Policy π_3

Dagger: Iteration N

Robot π_N drives



After many iterations
we are able to drive like a human!

DAgger (Dataset Aggregation)

Initialize with a random policy π_1 # Can be BC

Initialize empty data buffer $\mathcal{D} \leftarrow \{\}$

For $i = 1, \dots, N$

Execute policy π_i in the real world and collect data

$$\mathcal{D}_i = \{s_0, a_0, s_1, a_1, \dots\} \quad \# \text{ Also called a rollout}$$

Query the **expert** for the optimal action on **learner** states

$$\mathcal{D}_i = \{s_0, \pi^\star(s_0), s_1, \pi^\star(s_1), \dots\}$$

Aggregate data $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$

Train a new learner on this dataset $\pi_{i+1} \leftarrow \text{Train}(\mathcal{D})$

Select the best policy in $\pi_{1:N+1}$

The DAGGER Guarantee

DAGGER returns a policy π such that

$$\sum_{t=0}^{T-1} \mathbb{E}_{s_t \sim d_t^\pi} [\ell(s_t, \pi(s_t))] \leq O(\epsilon T)$$

Today's class

- ☑ What is DAGGER?
- ☐ DAGGER in the real world
- ☐ Why *aggregate* data?

Many cool applications of DAGGER in robotics



Lee et al, Learning quadrupedal locomotion over challenging terrain (2020)



Chen et al Learning by Cheating(2020)



Choudhury et al, Data Driven Planning via Imitation Learning (2018)



Pan et al Imitation learning for agile autonomous driving (2019)

How do we actually apply DAGGER in practice?

Asking a *human* expert to label every state
the robot visits is hard

Option 1: Extend DAGGER to different degrees of human feedback

Can we extend DAGGER to handle easier forms of human feedback preferences, interventions, etc?

Yes (*Future lectures!)

Option 2: Use an algorithmic oracle

What if we had a powerful algorithm
that we can run in train time
but not at test time?



Learning quadrupedal locomotion over challenging terrain

Joonho Lee¹, Jemin Hwangbo^{1,2†}, Lorenz Welhausen¹,
Vladlen Koltun³, Marco Hutter¹

¹ Robotic Systems Lab, ETH Zurich

² Robotics & Artificial Intelligence Lab, KAIST

³ Intelligent Systems Lab, Intel

†Substantial part of the work was carried out during his stay at 1

ETH zürich



Quiz: Is the teacher policy realizable?

When poll is active respond at PolleEv.com/sc2582

Send **sc2582** to **22333**



Today's class

- What is DAGGER?
- DAGGER in the real world
- Why *aggregate* data?

But why does
aggregating data work?



Why aggregate data?

Initialize with a random policy π_1

Initialize empty data buffer $\mathcal{D} \leftarrow \{\}$

For $i = 1, \dots, N$

Execute policy π_i in the real world and collect data

$$\mathcal{D}_i = \{s_0, a_0, s_1, a_1, \dots\}$$

Query the expert for the optimal action on learner states

$$\mathcal{D}_i = \{s_0, \pi^\star(s_0), s_1, \pi^\star(s_1), \dots\}$$

Aggregate data $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$

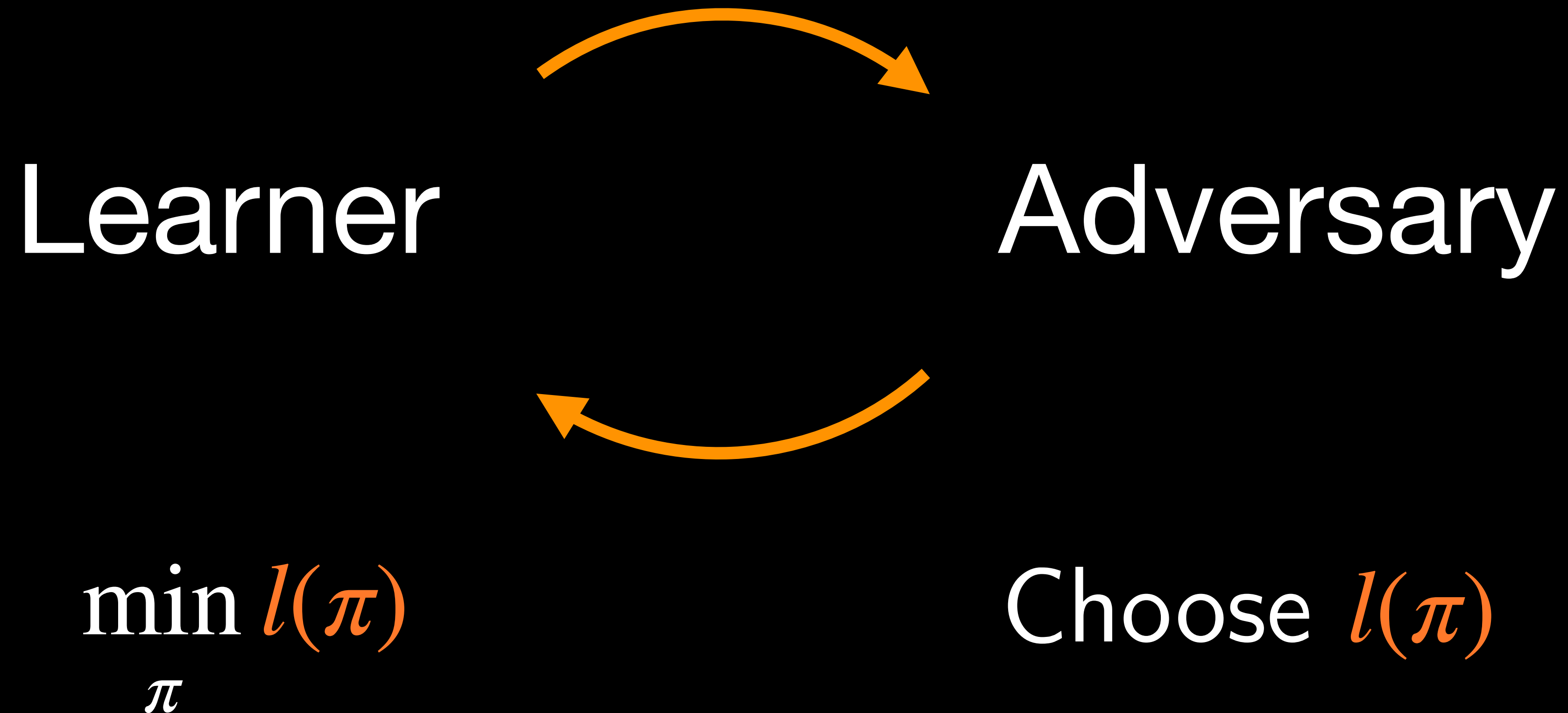
Train a new learner on this dataset $\pi_{i+1} \leftarrow \text{Train}(\mathcal{D})$

Select the best policy in $\pi_{1:N+1}$

From
Imitation Learning
to
Interactive
No-Regret Learning



Interactive Learning



Interactive Learning

Learner

Adversary

Initialize policy

π_1 [policy]

Chooses loss

$l_1(\cdot)$ [loss]

Update policy

π_2

Chooses loss

$l_2(\cdot)$

⋮

⋮

What is the best that I can do in such an adversarial setting?

From
Imitation Learning
to
Interactive
No-Regret Learning



How do we design algorithms that are no-regret?

$$\text{Regret} = \sum_{i=1}^N l_i(\pi_i) - \min_{\pi^*} \sum_{i=1}^N l_i(\pi^*)$$

(Learner)

(Best in
hindsight)

FOLLOW THE LEADER!



At every round t , choose the **best policy in hindsight**

$$\pi_i = \arg \min_{\pi} \sum_{j=1}^{i-1} l_j(\pi)$$

(lowest total loss)

$$\sum l_i$$

$$l_1$$

Policy 1

--



1.0

Policy 2

--



0.2

Policy 3

--



0.5

Avg. Regret: --

$\sum l_i$ l_1 l_2

Policy 1

1.0



1.0

0.5

Policy 2

0.2



0.2

0.5

Policy 3

0.5



0.5

0.2

Avg. Regret: 0.80

$$\sum l_i$$

$$l_1$$

$$l_2$$

$$l_3$$

1.5

Policy 1



1.0

0.5

0.5

0.7

Policy 2



0.2

0.5

1.0

0.7

Policy 3



0.5

0.2

0.2

Avg. Regret: **0.40**

$$\sum l_i$$

$$l_1$$

$$l_2$$

$$l_3$$

$$l_4$$

2.0

Policy 1



1.0

0.5

0.5

1.0

1.7

Policy 2



0.2

0.5

1.0

0.2

0.9

Policy 3



0.5

0.2

0.2

0.5

Avg. Regret: **0.53**

$$\sum l_i$$

$$l_1$$

$$l_2$$

$$l_3$$

$$l_4$$

$$l_5$$

Policy 1

3.0



1.0

0.5

0.5

1.0

0.5

Policy 2

1.9



0.2

0.5

1.0

0.2

1.0

Policy 3

1.4



0.5

0.2

0.2

0.5

0.2

Avg. Regret: 0.40

$$\sum l_i$$

$$l_1$$

$$l_2$$

$$l_3$$

$$l_4$$

$$l_5$$

$$l_6$$

Policy 1

3.5



1.0

0.5

0.5

1.0

0.5

1.0

Policy 2

2.9



0.2

0.5

1.0

0.2

1.0

0.5

Policy 3

1.6



0.5

0.2

0.2

0.5

0.2

0.2

Avg. Regret: 0.32

$$\sum l_i$$

$$l_1$$

$$l_2$$

$$l_3$$

$$l_4$$

$$l_5$$

$$l_6$$

Policy 1

4.5



1.0

0.5

0.5

1.0

0.5

1.0

Policy 2

3.4



0.2

0.5

1.0

0.2

1.0

0.5

Policy 3

1.8



0.5

0.2

0.2

0.5

0.2

0.2

Avg. Regret: 0.26

Is FTL no-regret?

FTL is no-regret if

1. We are in the continuous setting
2. Loss is strongly convex

Back to the
proof!

