

# Inverse Reinforcement Learning: From Maximum Margin to Maximum Entropy

Sanjiban Choudhury



Cornell Bowers CIS  
**Computer Science**

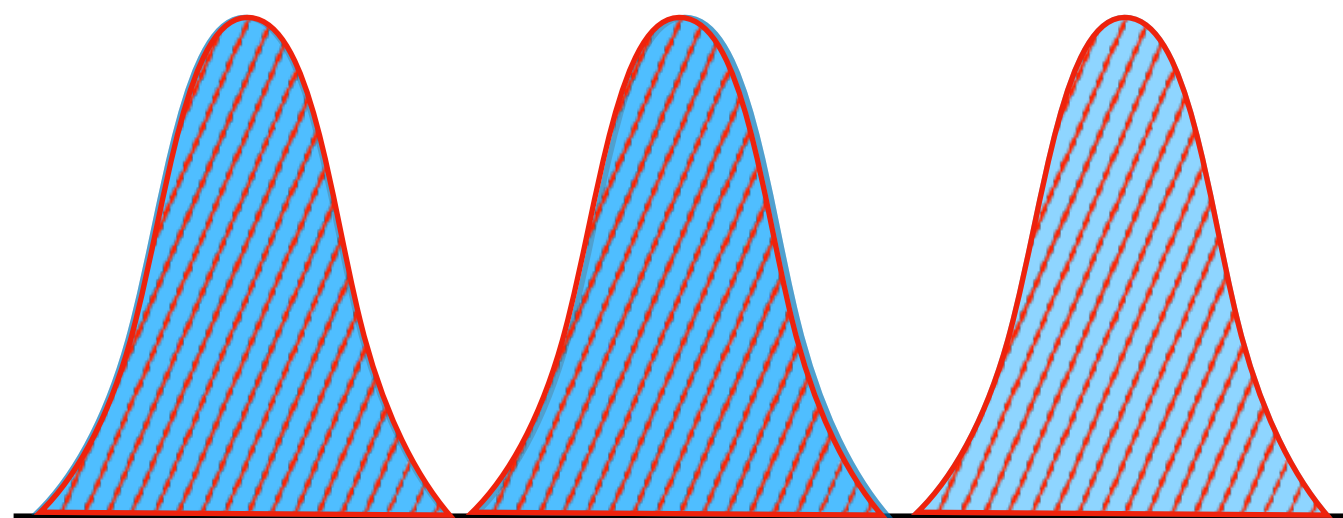
Easy



Expert is **realizable**

$$\pi^E \in \Pi$$

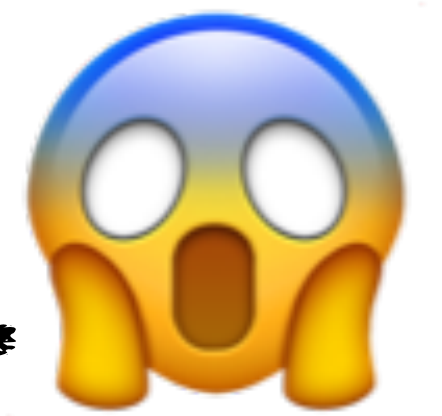
As  $N \rightarrow \infty$ , drive down  
 $\epsilon = 0$  (or Bayes error)



Nothing special.

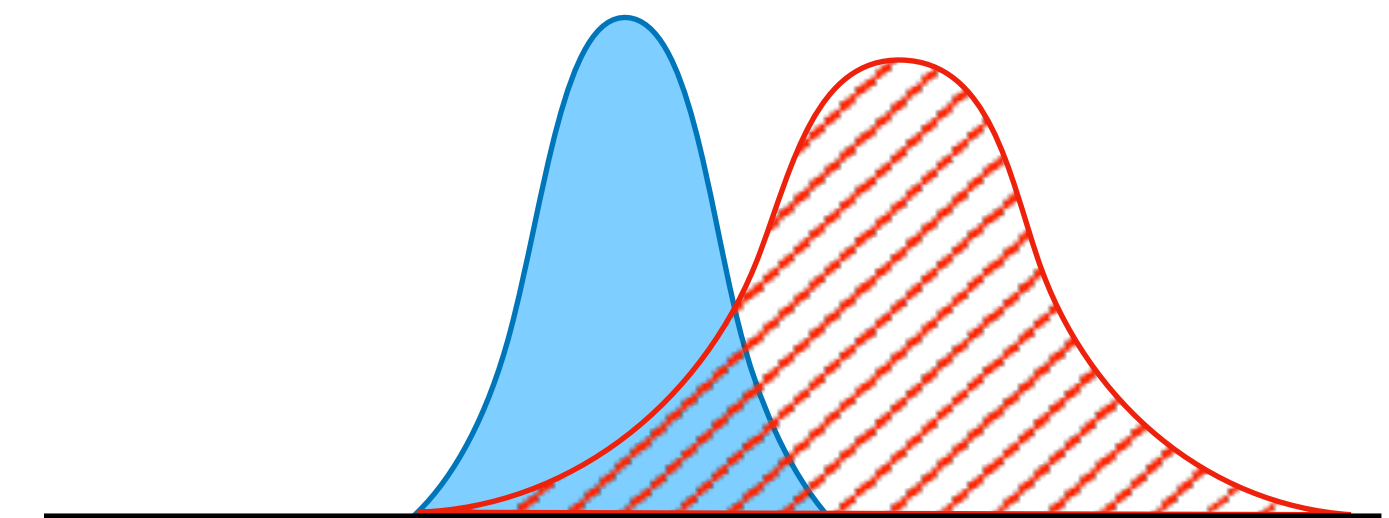
Collect lots of data and  
do Behavior Cloning

Hard



Non-realizable expert +  
limited expert support

Even as  $N \rightarrow \infty$ ,  
behavior cloning  $O(\epsilon T^2)$

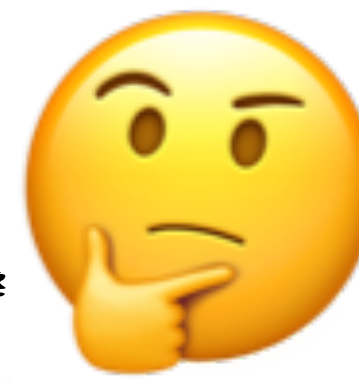


Requires **interactive** expert  
(DAGGER/EIL) to provide  
labels  $\Rightarrow O(\epsilon T)$

Easy



Medium



Hard



Expert is **realizable**

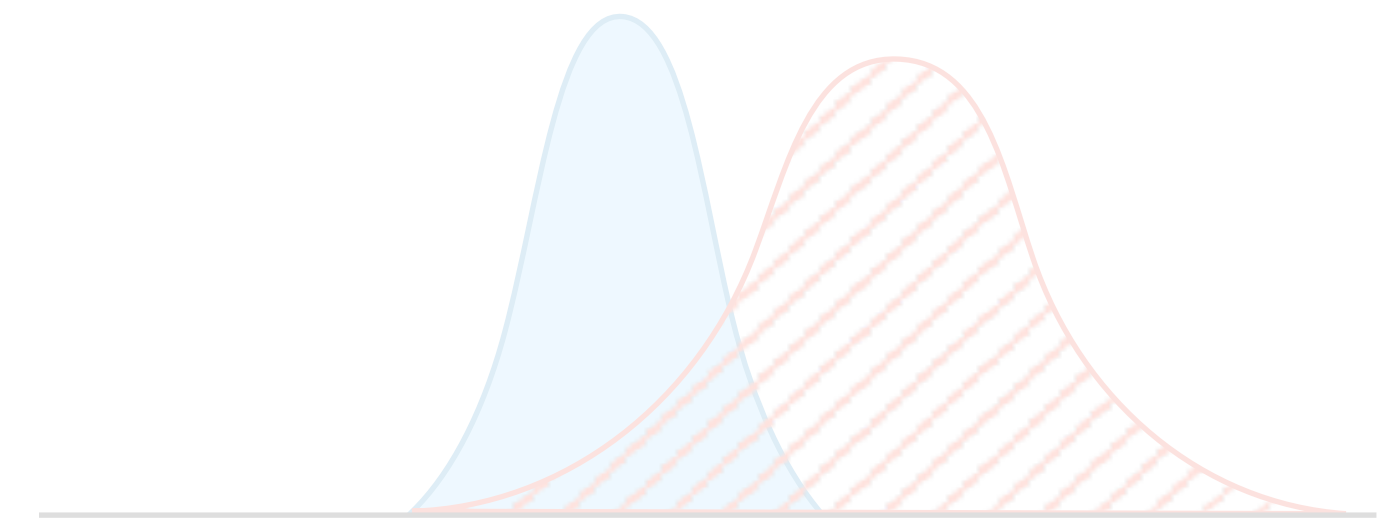
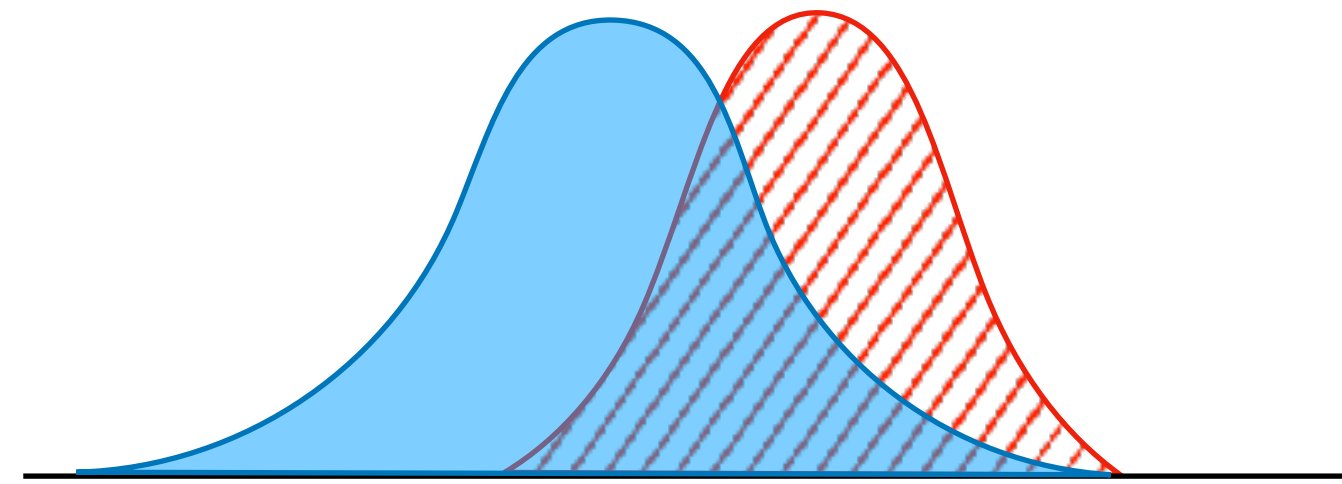
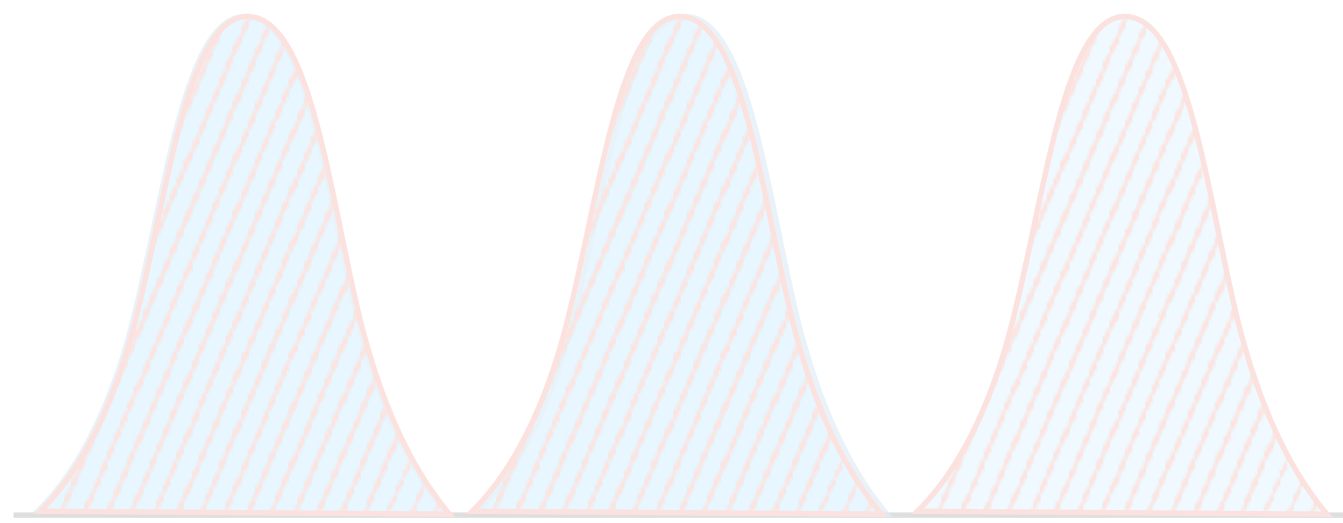
$$\pi^E \in \Pi$$

Non-realizable expert  
but **full expert support**

Non-realizable expert +  
limited expert support

As  $N \rightarrow \infty$ , drive down  
 $\epsilon = 0$  (or Bayes error)

Even as  $N \rightarrow \infty$ ,  
behavior cloning  $O(\epsilon T^2)$



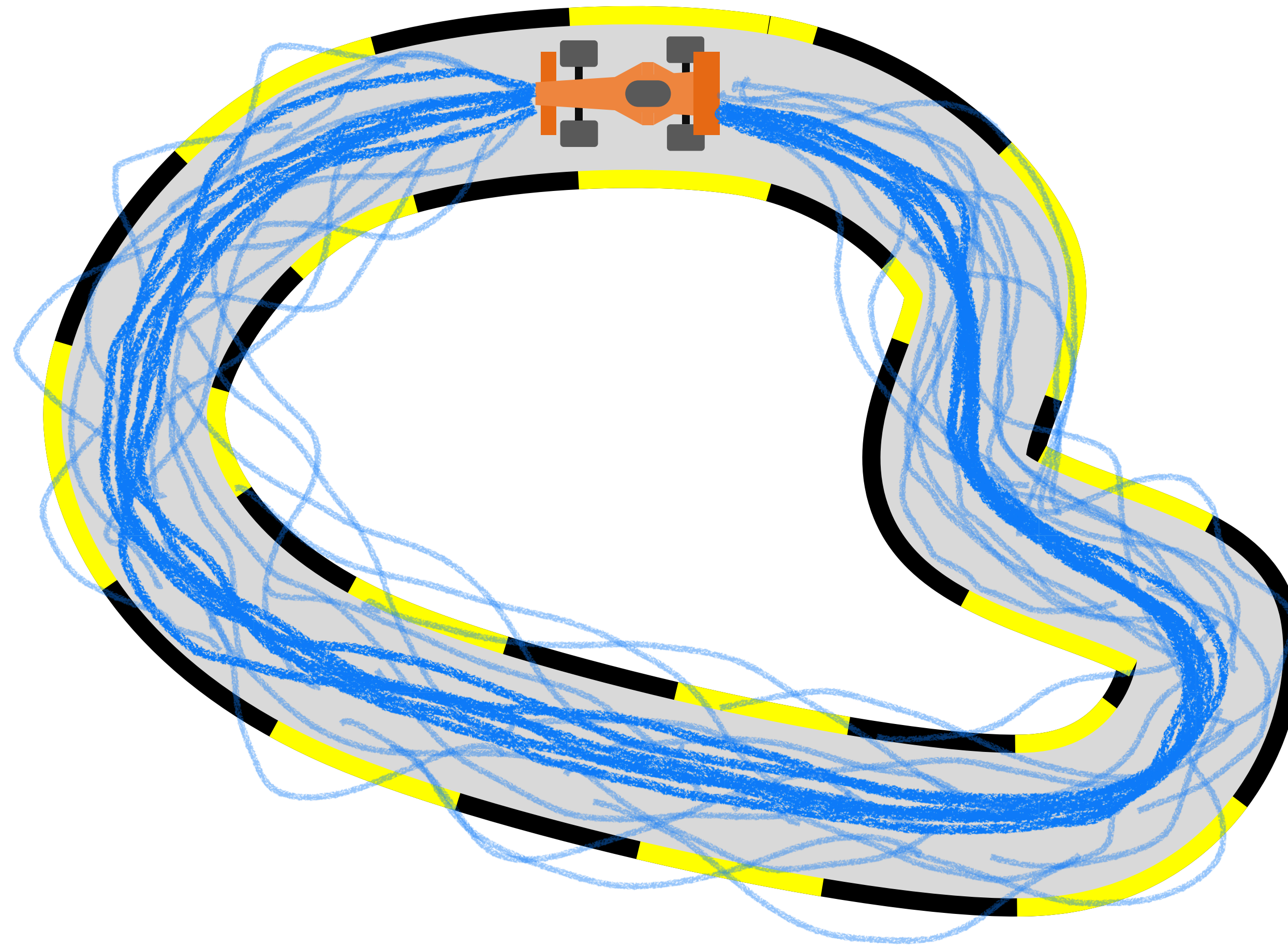
Nothing special.

Collect lots of data and  
do Behavior Cloning



Requires **interactive** expert  
(DAGGER/EIL) to provide  
labels  $\Rightarrow O(\epsilon T)$

# Expert demonstrations have full coverage



.. but expert runs away after demonstrations

So expert data has  
full coverage ...

.. why don't we just do  
Behavior Cloning?



Activity!



# Think-Pair-Share!

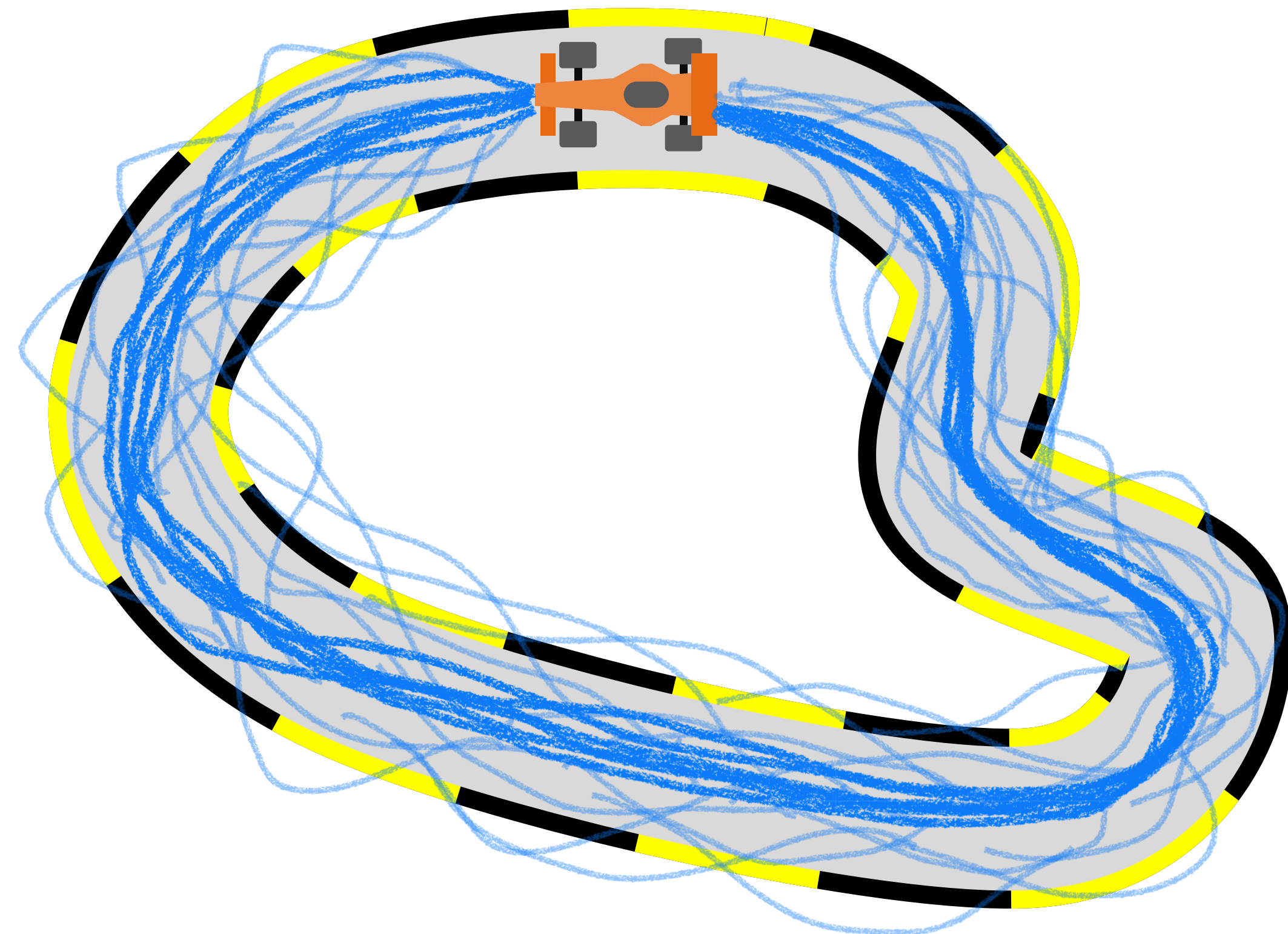


*To BC, or not to BC!*

Think (30 sec): Will BC work?  $O(\epsilon T)$  or  $O(\epsilon T^2)$  ?  
Make the argument!

Pair: Find a partner

Share (45 sec): Partners exchange ideas





BC results in compounding  
errors

We don't have an interactive  
expert





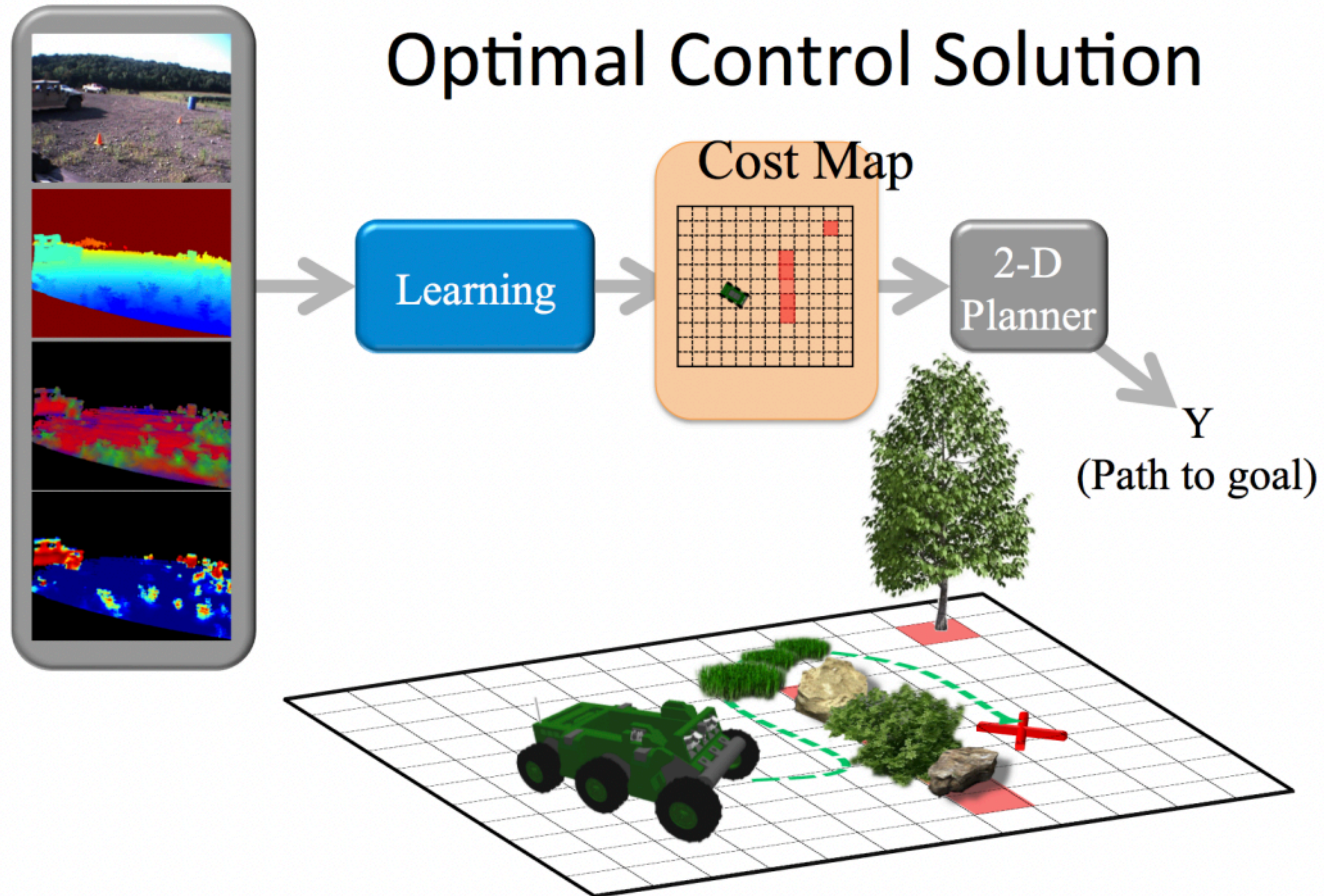
What if we knew our MDP  
(except the cost)?

Or what if we had an  
**interactive** simulator?

# CRUSHER robot from CMU



# Can we learn a cost function for CRUSHER navigation?



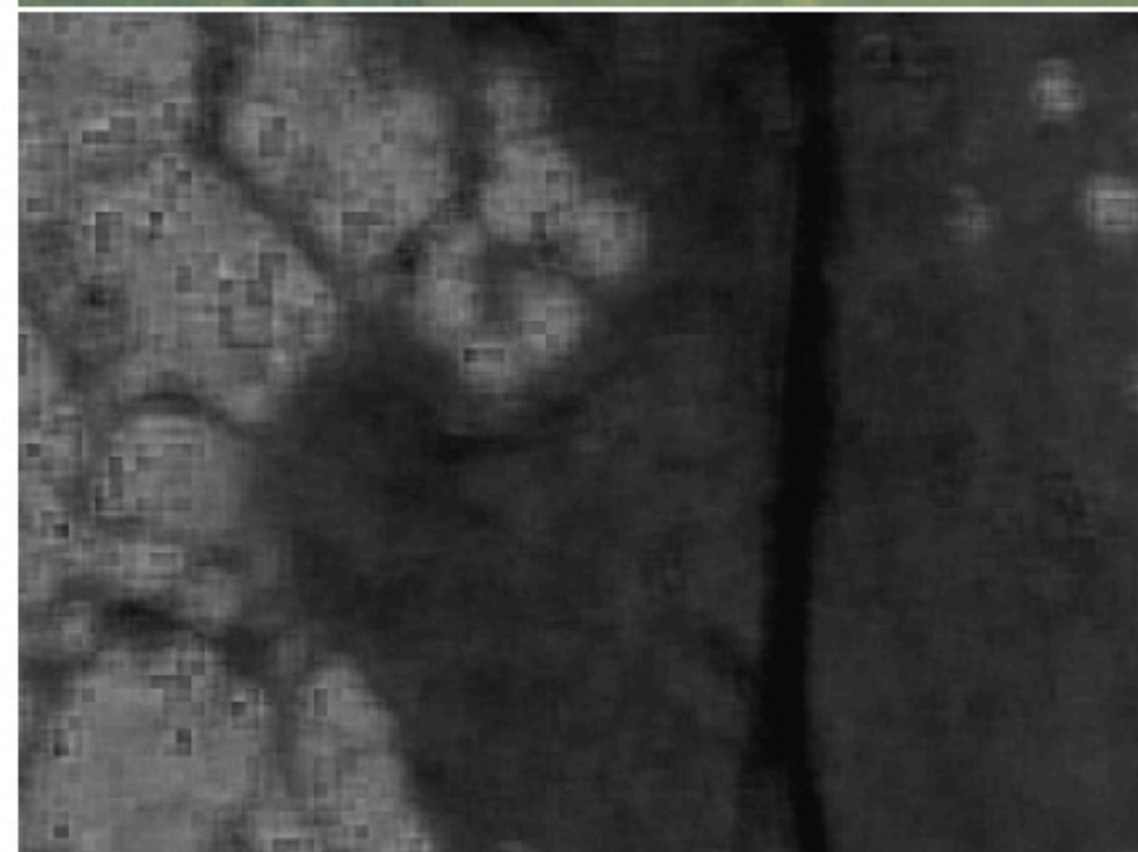
Let's  
formalize!



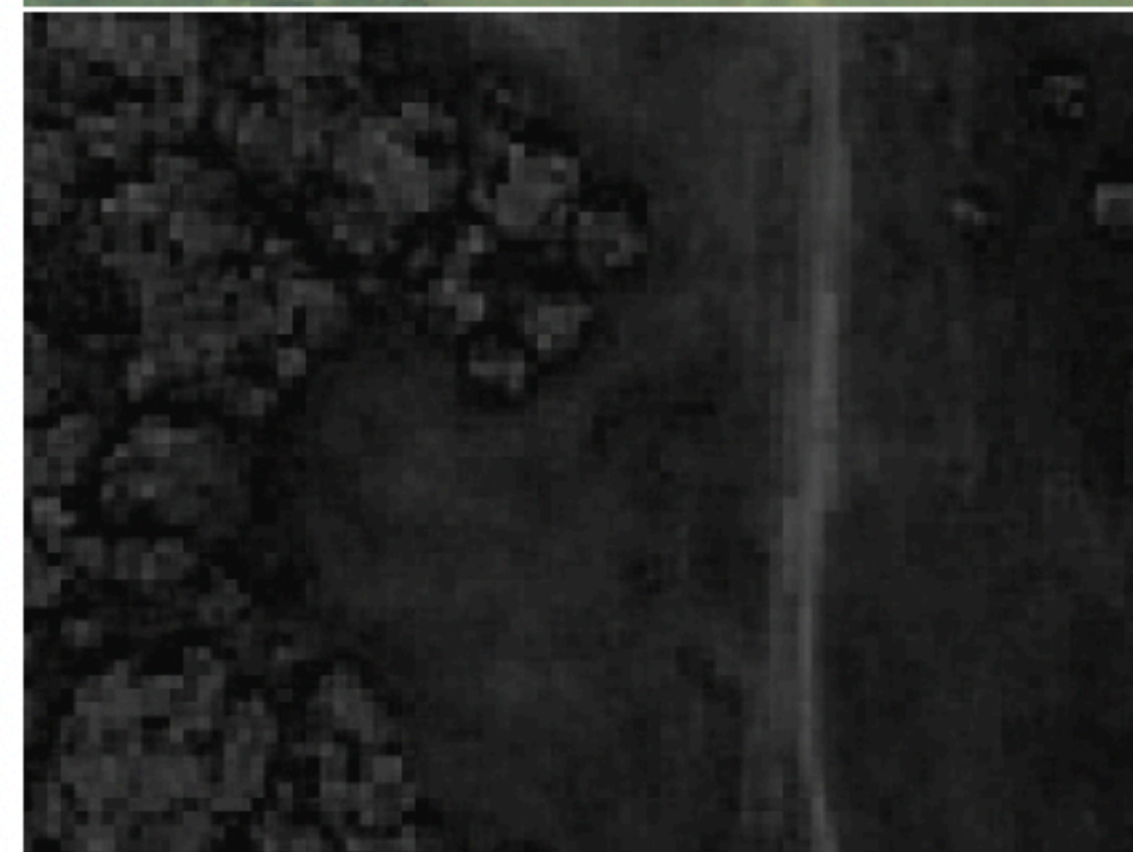
# Learning to Search (LEARCH)



Min distance



Stay on roads



Stay near trees



## Learning to Search: Functional Gradient Techniques for Imitation Learning

Nathan D. Ratliff  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213  
ndr@ri.cmu.edu

David Silver  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213  
dsilver@ri.cmu.edu

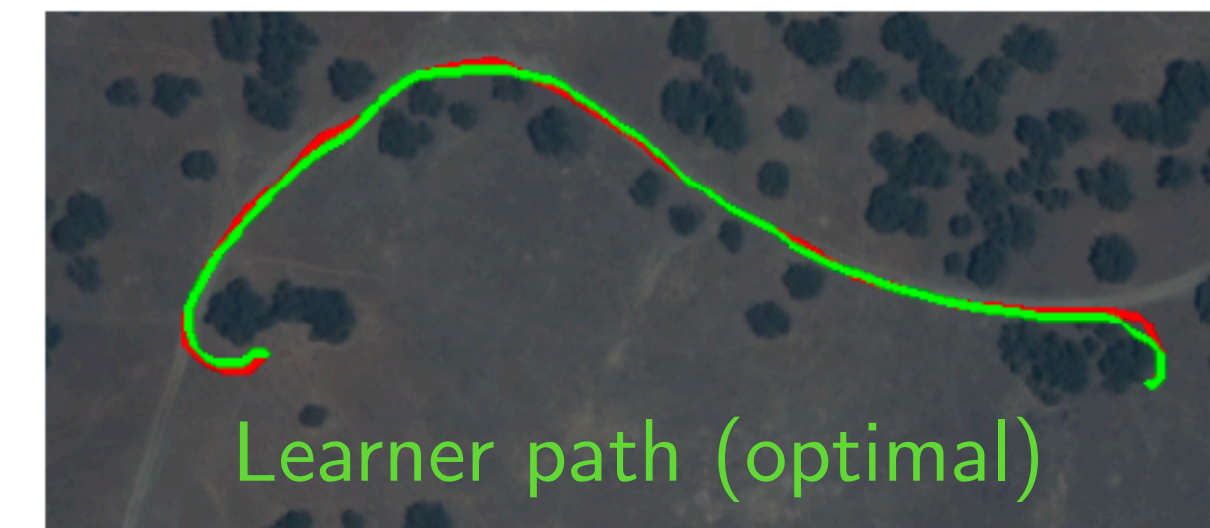
J. Andrew Bagnell  
Robotics Institute and Machine Learning  
Carnegie Mellon University  
Pittsburgh, PA 15213  
dbagnell@ri.cmu.edu

# Learning to Search (LEARCH)

Human demonstration



Human demonstration



Given dataset:  $\left\{ \underset{\text{(Human demo)}}{\xi_i^h}, \underset{\text{(Map)}}{\phi_i} \right\}_{i=1}^N$

Solve for cost  $C_{\theta}(\xi)$

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \left( \underset{\text{Human Cost}}{C_{\theta}(\xi_i^h, \phi_i)} - \min_{\xi} \left[ \underset{\text{Learner Cost}}{C_{\theta}(\xi, \phi_i)} - \underset{\text{(Margin)}}{\gamma(\xi, \xi^h)} \right] \right) + R(\theta)$$

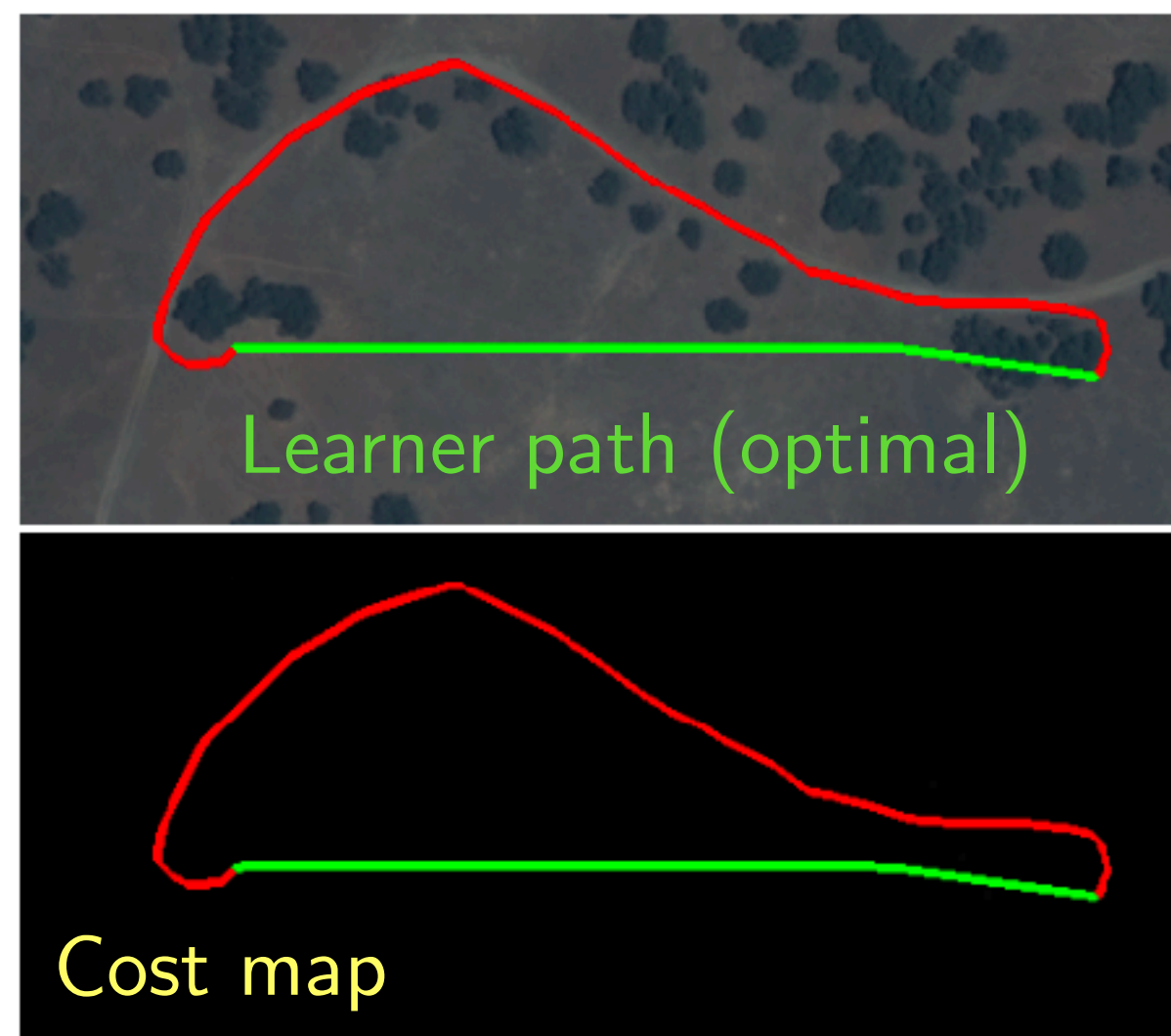
Human  
Cost

Learner  
Cost

Regularizer

# Learning to Search (LEARCH)

Human demonstration



for  $i = 1, \dots, N$

# Loop over datapoints

$$\xi_i^* = \min_{\xi} [C_{\theta}(\xi, \phi_i) - \gamma(\xi, \xi^h)]$$

# Call planner!

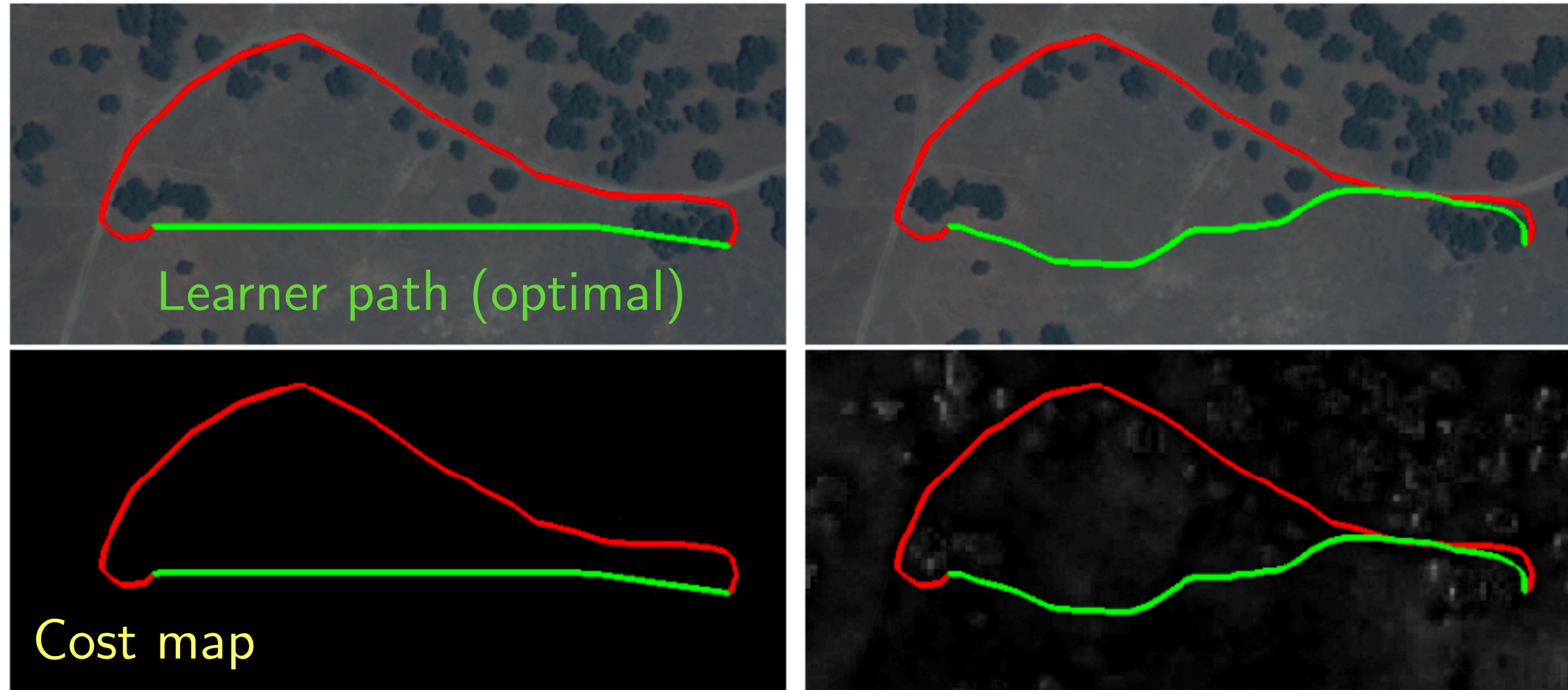
$$\theta^+ = \theta - \eta \left[ \nabla_{\theta} C_{\theta}(\xi_i^h, \phi_i) - \nabla_{\theta} C_{\theta}(\xi_i^*, \phi_i) + \nabla_{\theta} R(\theta) \right]$$

(Push down human cost) (Push up planner cost)

# Update cost

# Learning to Search

Human demonstration



for  $i = 1, \dots, N$

# Loop over datapoints

$$\xi_i^* = \min_{\xi} [C_{\theta}(\xi, \phi_i) - \gamma(\xi, \xi^h)]$$

# Call planner!

$$\theta^+ = \theta - \eta \left[ \nabla_{\theta} C_{\theta}(\xi_i^h, \phi_i) - \nabla_{\theta} C_{\theta}(\xi_i^*, \phi_i) + \nabla_{\theta} R(\theta) \right]$$

(Push down human cost)

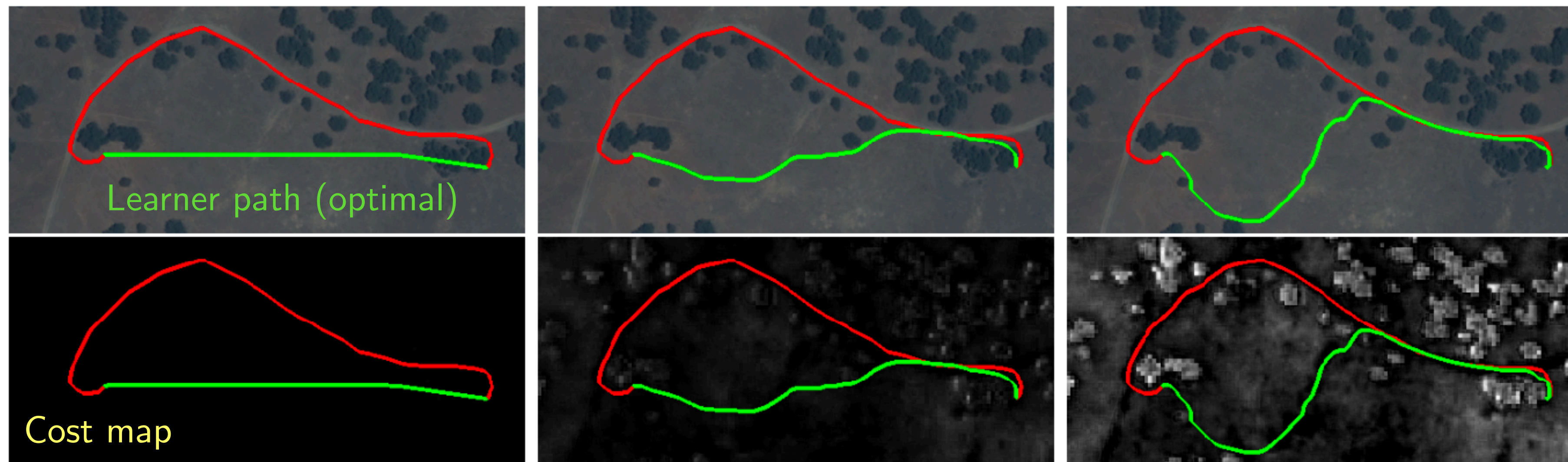
(Push up planner cost)

# Update cost



# Learning to Search (LEARCH)

Human demonstration



for  $i = 1, \dots, N$

# Loop over datapoints

$$\xi_i^* = \min_{\xi} [C_{\theta}(\xi, \phi_i) - \gamma(\xi, \xi^h)]$$

# Call planner!

$$\theta^+ = \theta - \eta \left[ \nabla_{\theta} C_{\theta}(\xi_i^h, \phi_i) - \nabla_{\theta} C_{\theta}(\xi_i^*, \phi_i) + \nabla_{\theta} R(\theta) \right]$$

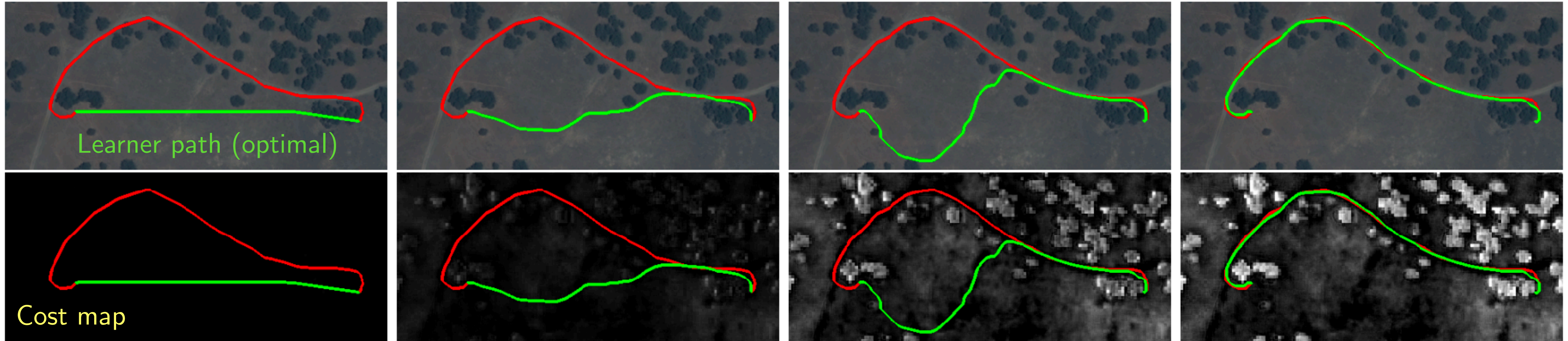
(Push down human cost)

(Push up planner cost)

# Update cost

# Learning to Search (LEARCH)

Human demonstration



for  $i = 1, \dots, N$

# Loop over datapoints

$$\xi_i^* = \min_{\xi} [C_{\theta}(\xi, \phi_i) - \gamma(\xi, \xi^h)]$$

# Call planner!

$$\theta^+ = \theta - \eta [ \underbrace{\nabla_{\theta} C_{\theta}(\xi_i^h, \phi_i)}_{\text{(Push down human cost)}} - \underbrace{\nabla_{\theta} C_{\theta}(\xi_i^*, \phi_i)}_{\text{(Push up planner cost)}} + \nabla_{\theta} R(\theta) ]$$

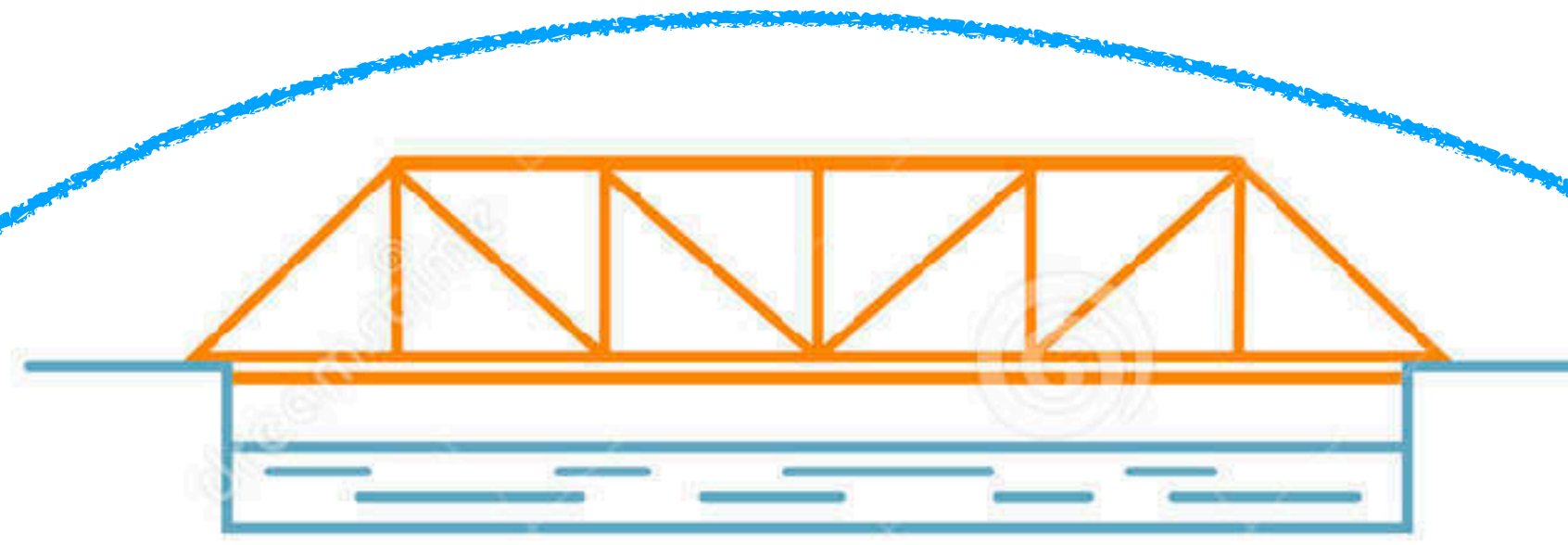
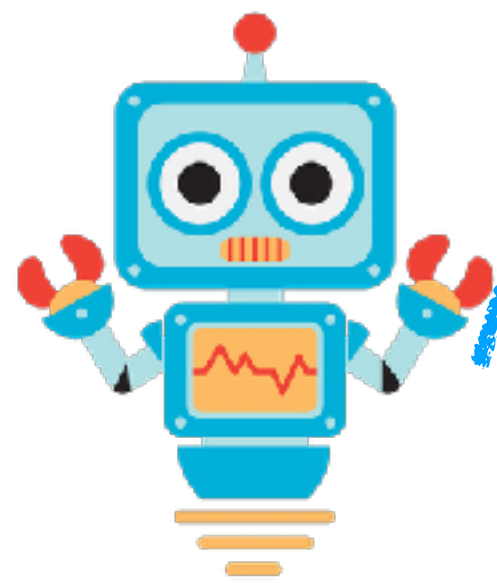
# Update cost

What happens when the expert is stochastic / noisy / suboptimal?

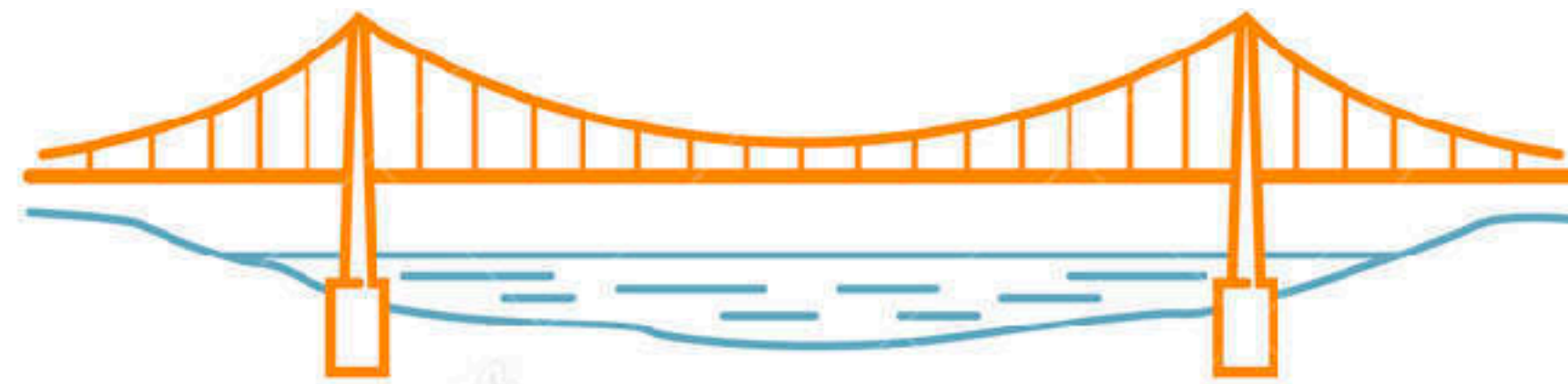


# Learning which bridge to cross

Demonstrations  
always pick  
Bridge 1



$$f = [1 \ 0 \ 0]$$



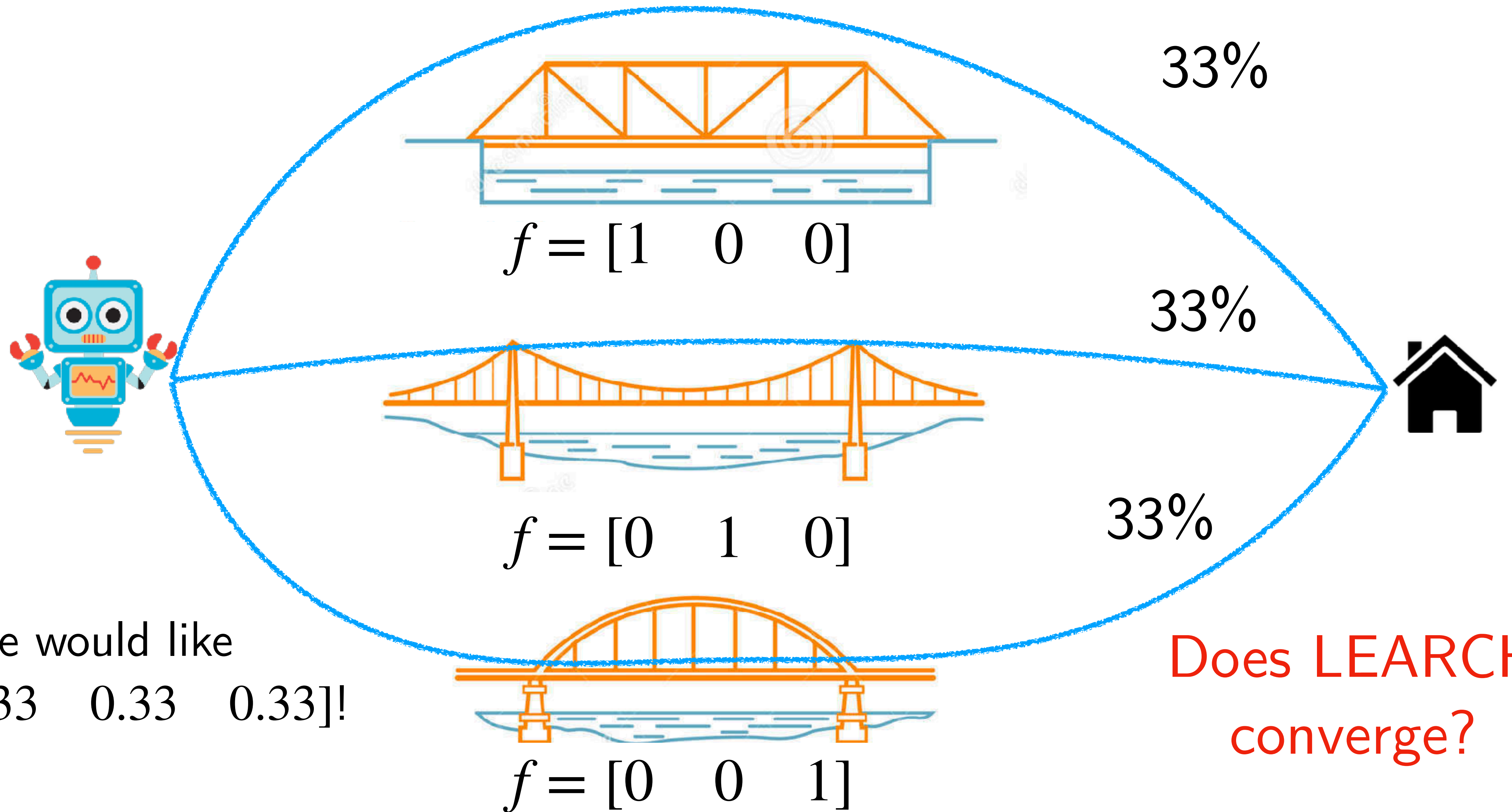
$$f = [0 \ 1 \ 0]$$



$$f = [0 \ 0 \ 1]$$

LEARCH converges to  
 $w = [1 \ 0 \ 0]!$

# Learning which bridge to cross





When the expert

is Suboptimal /

is Noisy /

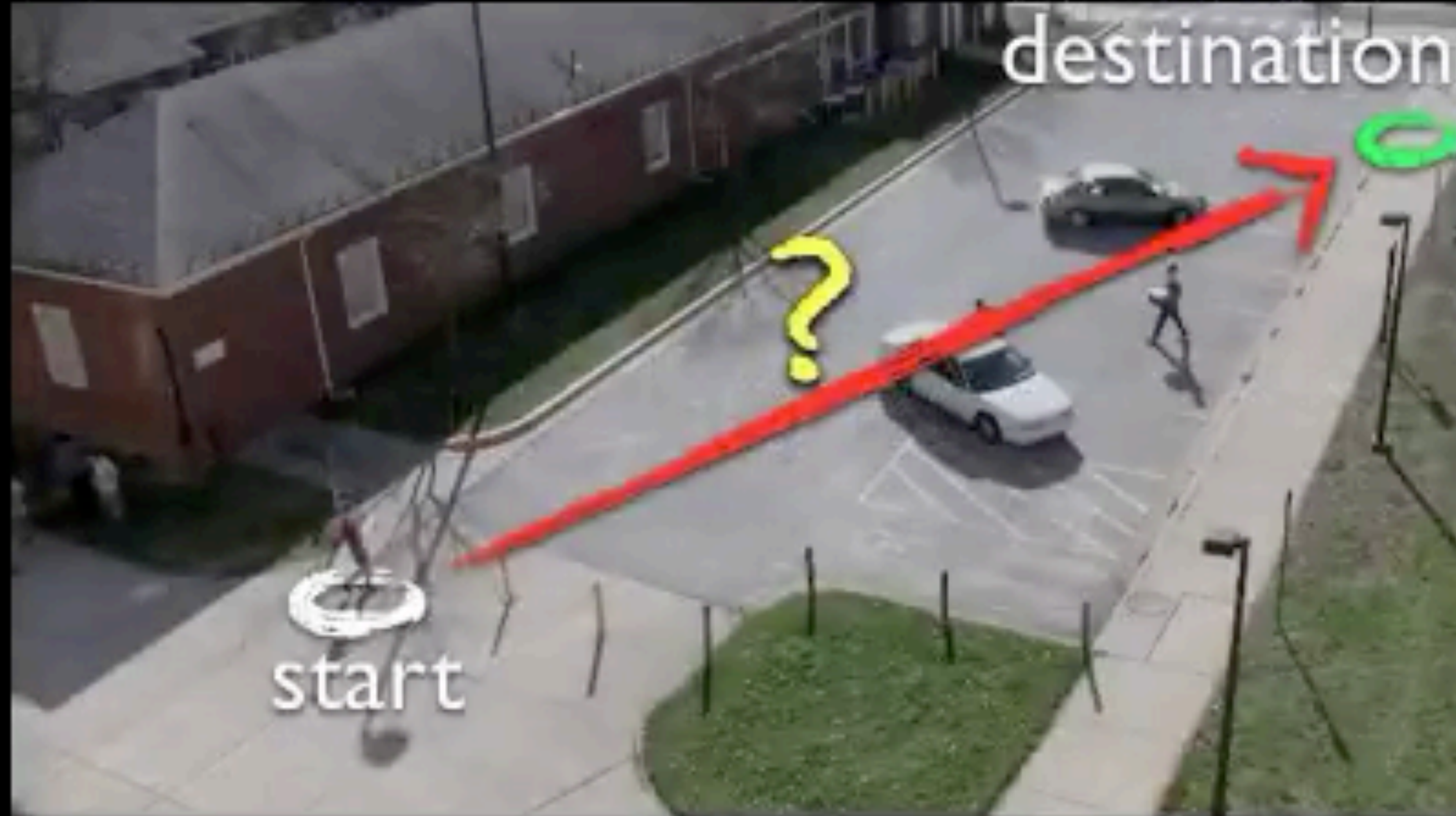
Has Privileged Information

LEARCH does NOT converge!!

# Activity Forecasting

Kris M. Kitani, Brian D. Ziebart, J. Andrew Bagnell, and Martial Hebert

Carnegie Mellon University, Pittsburgh, PA 15213 USA  
{kkitani,bziebart}@cs.cmu.edu, {dbagnell,hebert}@ri.cmu.edu





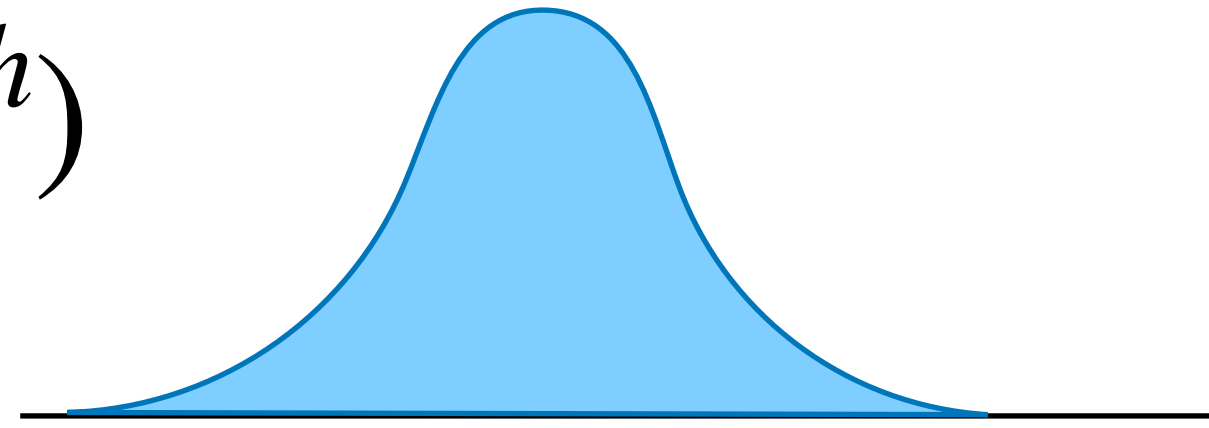
Expert demonstrations are  
coming from some (unknown)  
distribution ..

Can we learn this distribution?



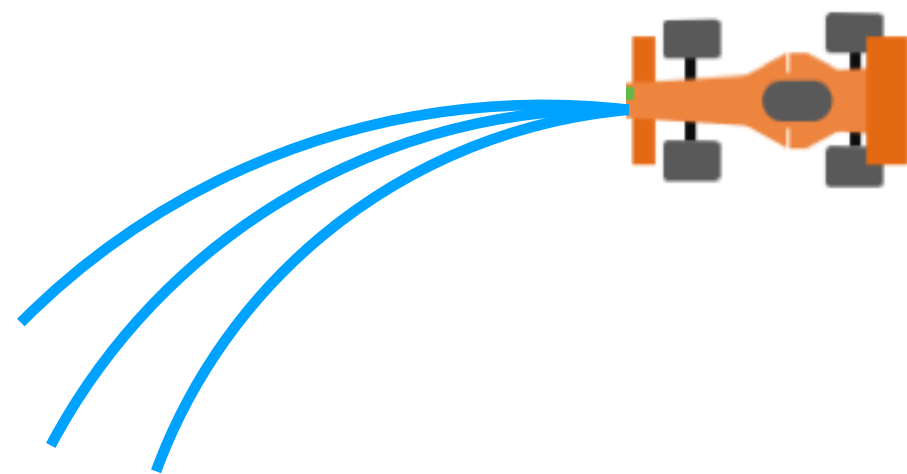
# The Distribution Matching Problem

$$P_{expert}(\xi^h)$$



(Unknown) expert distribution

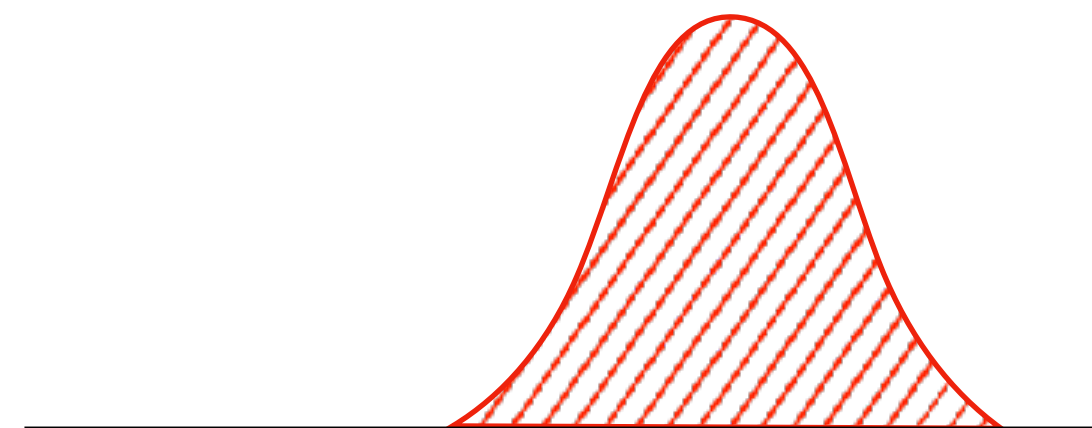
All we see are expert samples



What loss should we use?



$$P_{\theta}(\xi)$$



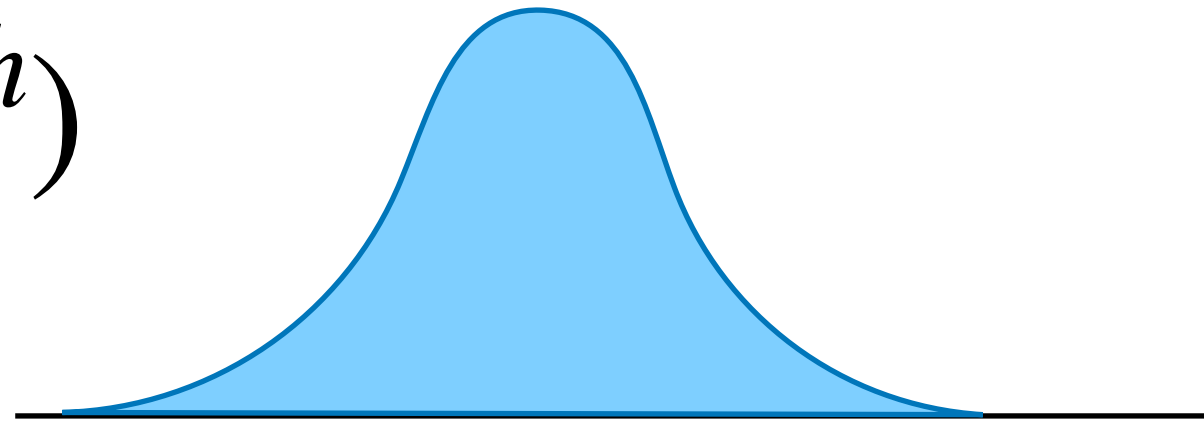
Learn distribution over trajectories

Learner can also generate samples



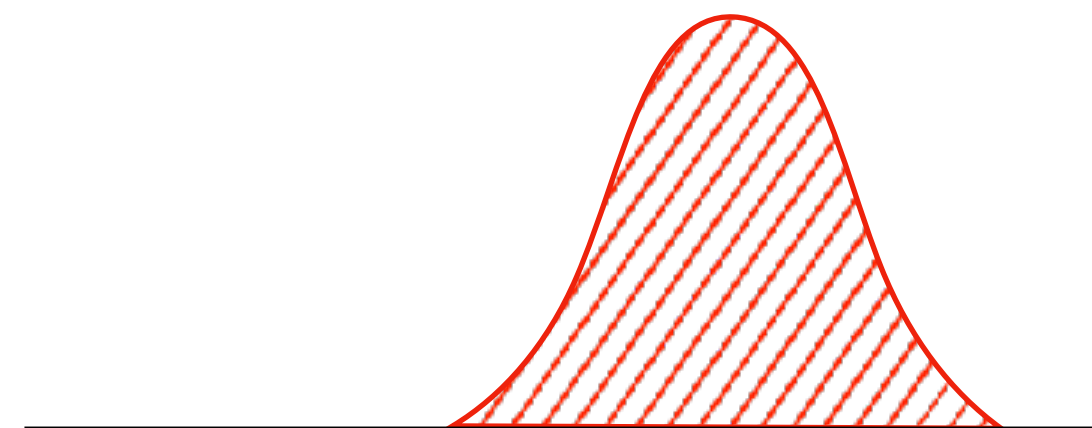
# Proposal: Match expected features!

$P_{expert}(\xi^h)$



(Unknown) expert distribution

$P_{\theta}(\xi)$



Learn distribution over trajectories

All we see are expert samples



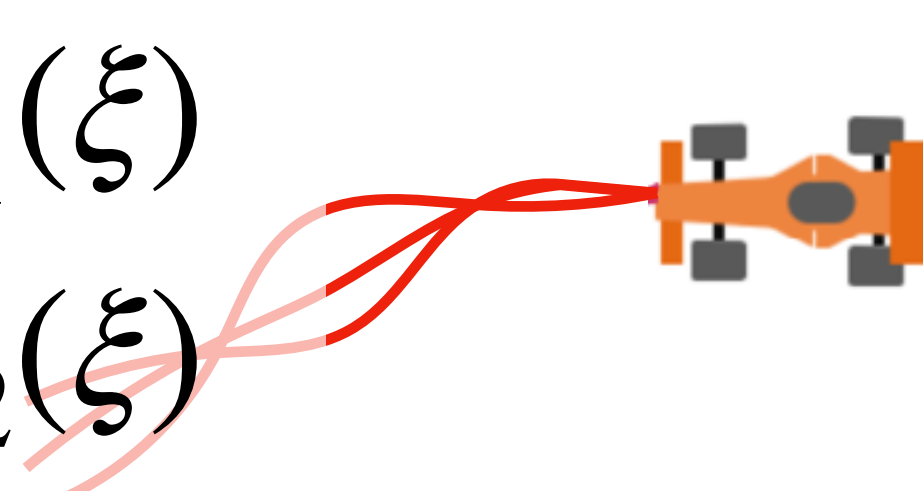
$$\mathbb{E}_{\xi^h \sim P_{expert}(\cdot)} f_1(\xi^h) = \mathbb{E}_{\xi \sim P_{\theta}(\cdot)} f_1(\xi)$$

$$\mathbb{E}_{\xi^h \sim P_{expert}(\cdot)} f_2(\xi^h) = \mathbb{E}_{\xi \sim P_{\theta}(\cdot)} f_2(\xi)$$

⋮

$$\mathbb{E}_{\xi^h \sim P_{expert}(\cdot)} f_k(\xi^h) = \mathbb{E}_{\xi \sim P_{\theta}(\cdot)} f_k(\xi)$$

Learner can also generate samples



Let's  
formalize!



# Maximum Entropy Inverse Optimal Control

## **Maximum Entropy Inverse Reinforcement Learning**

**Brian D. Ziebart, Andrew Maas, J.Andrew Bagnell, and Anind K. Dey**

School of Computer Science

Carnegie Mellon University

Pittsburgh, PA 15213

bziebart@cs.cmu.edu, amaas@andrew.cmu.edu, dbagnell@ri.cmu.edu, anind@cs.cmu.edu

# Maximum Entropy Inverse Optimal Control



## **LEO: Learning Energy-based Models in Factor Graph Optimization**

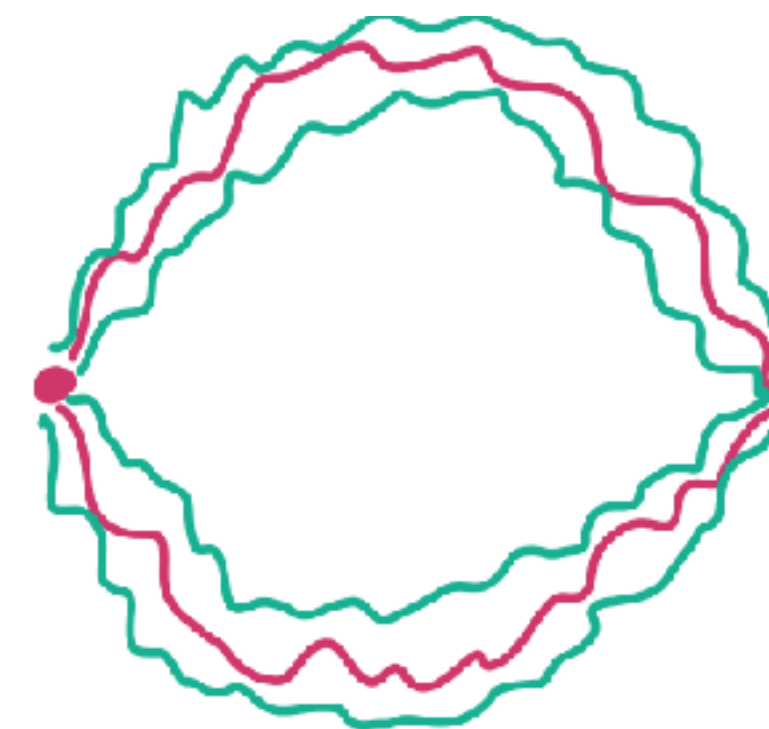
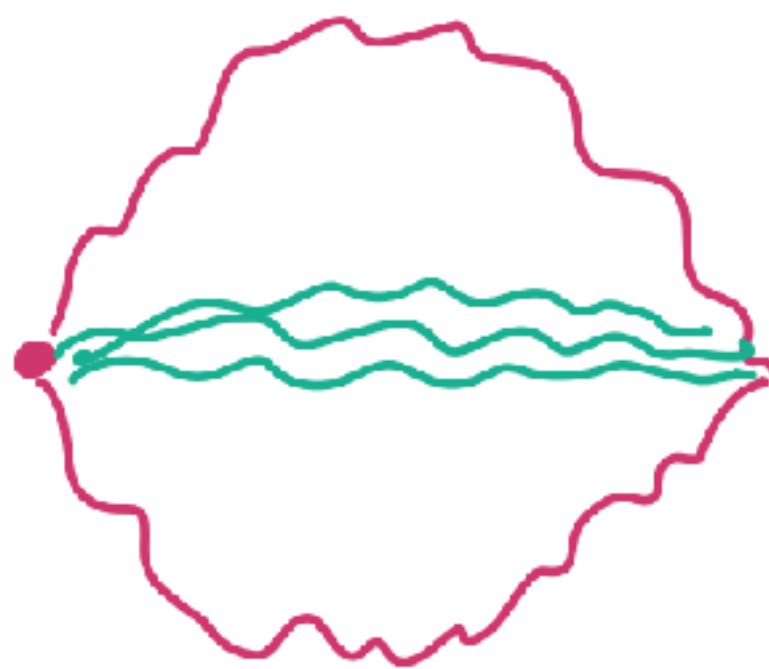
**Paloma Sodhi<sup>1,2</sup>, Eric Dexheimer<sup>1</sup>, Mustafa Mukadam<sup>2</sup>, Stuart Anderson<sup>2</sup>, Michael Kaess<sup>1</sup>**

<sup>1</sup>Carnegie Mellon University, <sup>2</sup>Facebook AI Research

# Maximum Entropy Inverse Optimal Control

Human demonstration

Learner traj



Given dataset:  $\left\{ \underset{\text{(Human demo)}}{\xi_i^h}, \underset{\text{(Map)}}{\phi_i} \right\}_{i=1}^N$

Solve for cost  $C_\theta(\xi)$

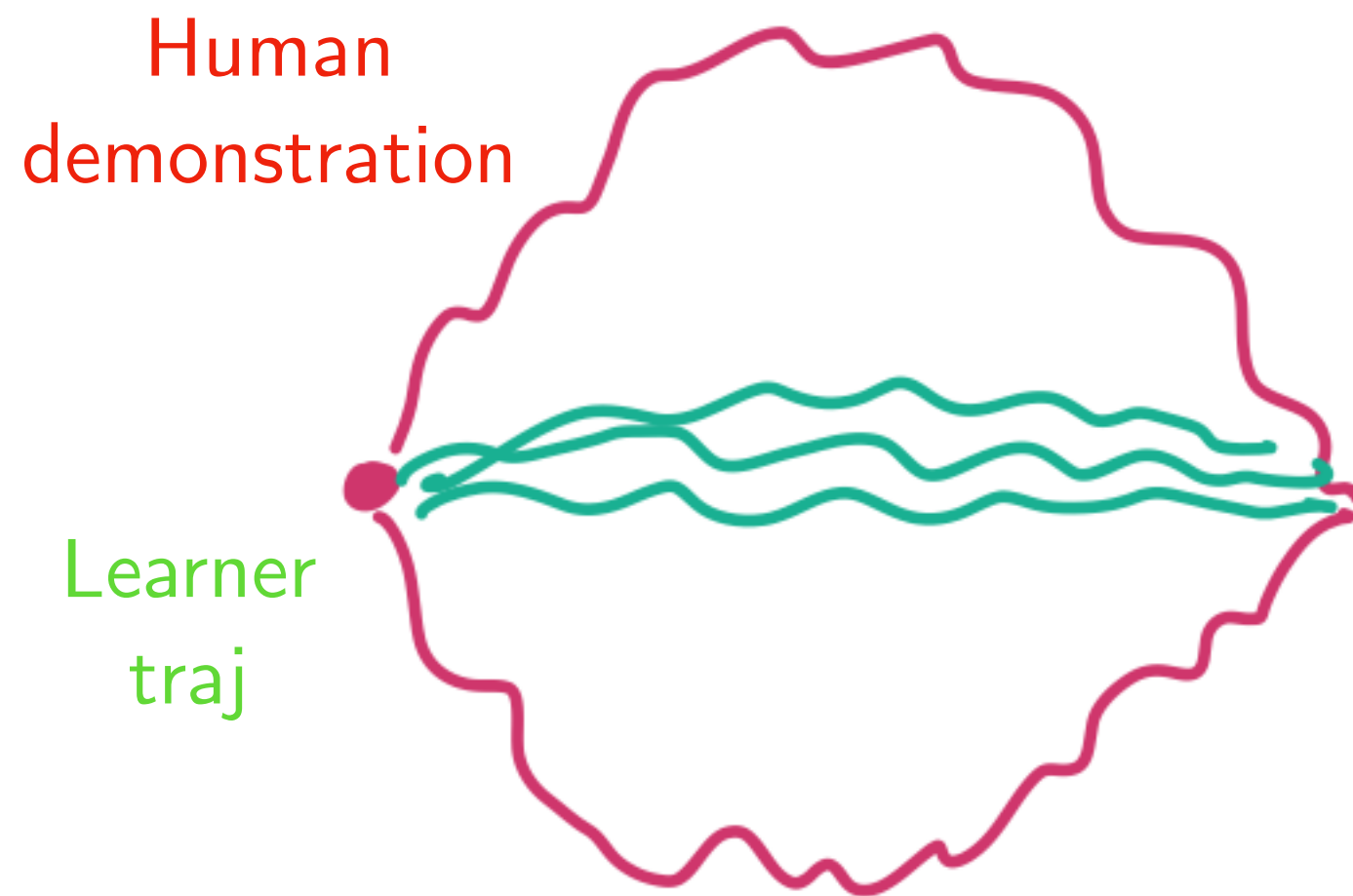
$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N -\log P_{\theta}(\xi_i^h | \phi_i)$$

*Max lik. of human traj*

$$P_{\theta}(\xi | \phi) = \frac{1}{Z(\theta, \phi)} \exp(-C_{\theta}(\xi, \phi))$$

*More costly traj, less likely*

# Maximum Entropy Inverse Optimal Control



for  $i = 1, \dots, N$

# Loop over datapoints

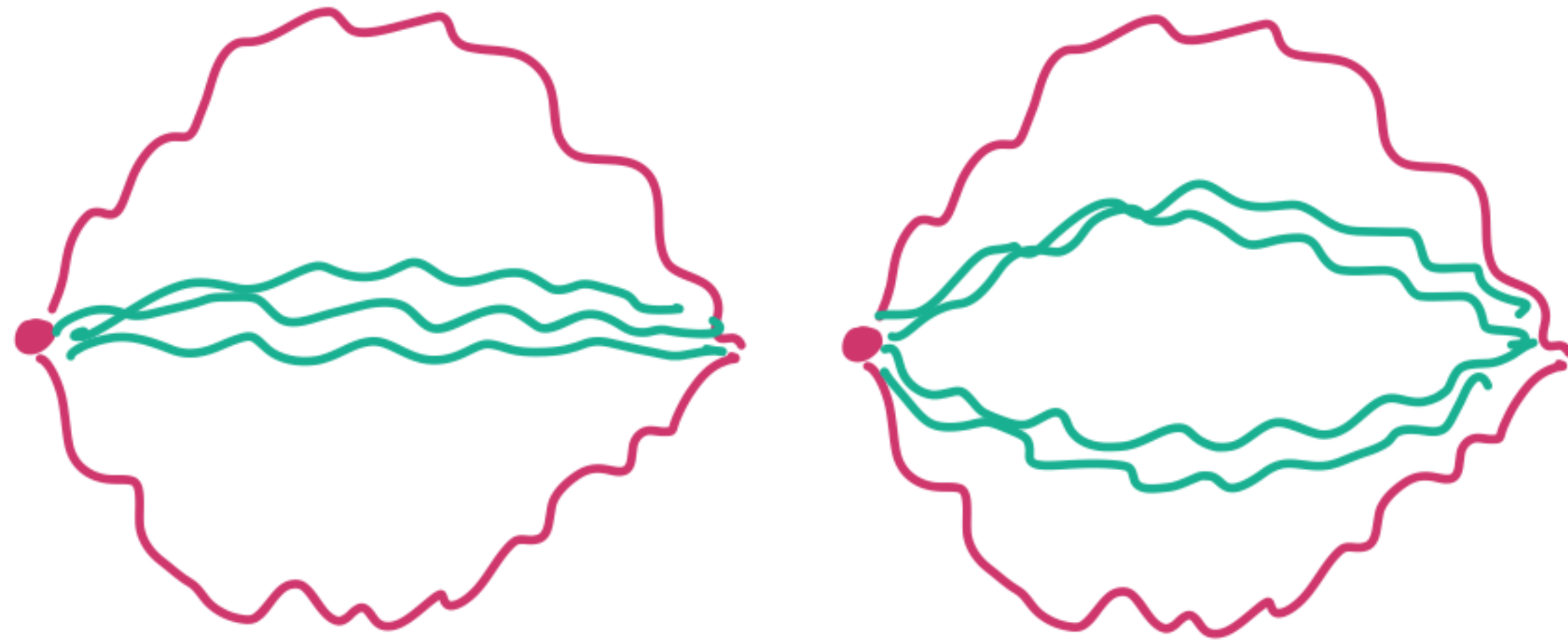
$$\xi_i \sim \frac{1}{Z} \exp(-C_\theta(\xi, \phi_i))$$

# Call planner!

$$\theta^+ = \theta - \eta \left[ \underbrace{\nabla_\theta C_\theta(\xi_i^h, \phi_i)}_{\text{(Push down human cost)}} - \underbrace{\nabla_\theta C_\theta(\xi_i, \phi_i)}_{\text{(Push up planner cost)}} \right]$$

# Update cost

# Maximum Entropy Inverse Optimal Control



for  $i = 1, \dots, N$

# Loop over datapoints

$$\xi_i \sim \frac{1}{Z} \exp(-C_\theta(\xi, \phi_i))$$

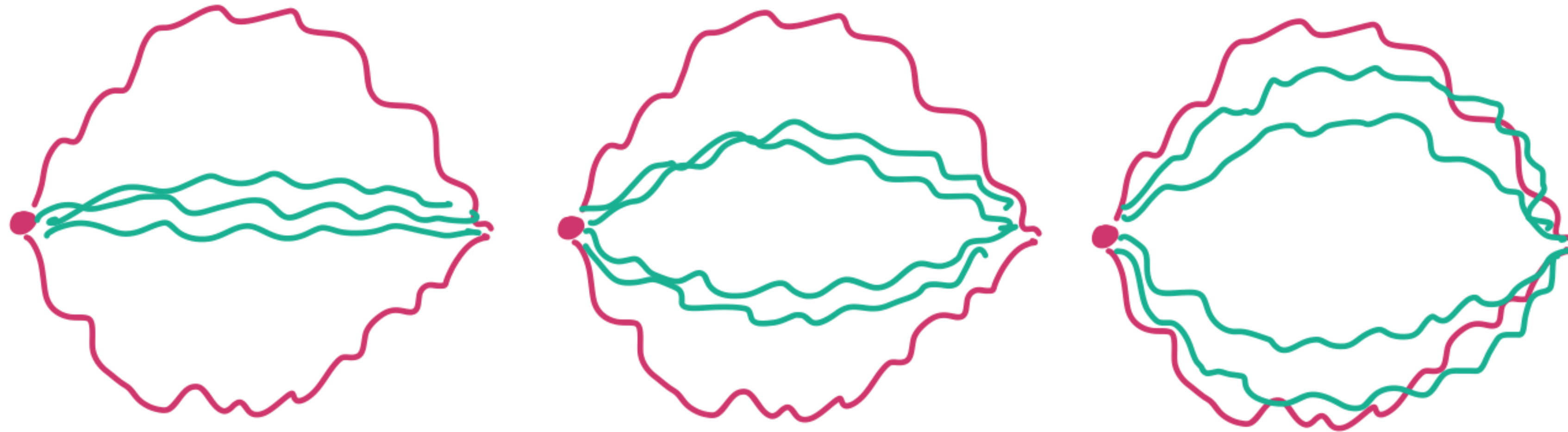
# Call planner!

$$\theta^+ = \theta - \eta \left[ \underbrace{\nabla_\theta C_\theta(\xi_i^h, \phi_i)}_{\text{(Push down human cost)}} - \underbrace{\nabla_\theta C_\theta(\xi_i, \phi_i)}_{\text{(Push up planner cost)}} \right]$$

# Update cost



# Maximum Entropy Inverse Optimal Control



for  $i = 1, \dots, N$

# Loop over datapoints

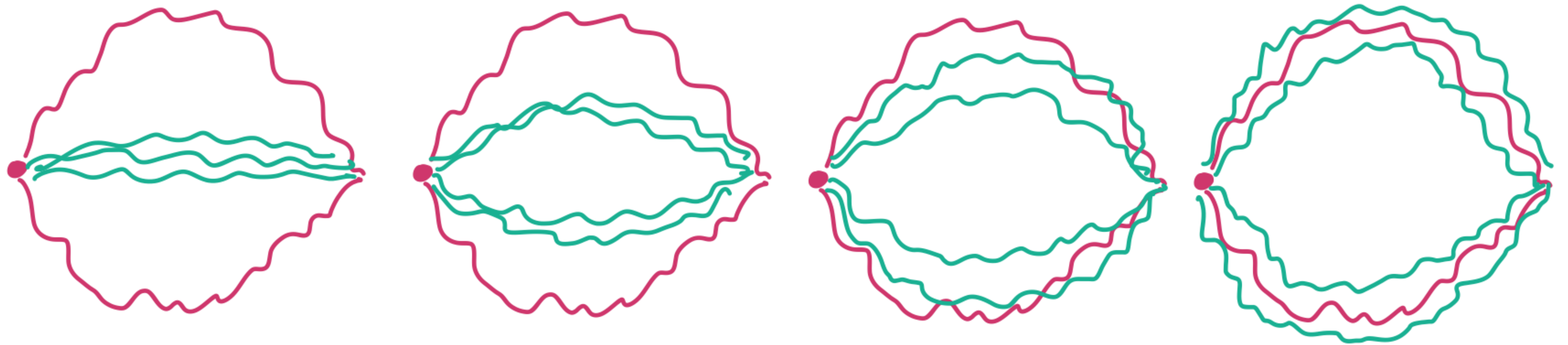
$$\xi_i \sim \frac{1}{Z} \exp(-C_\theta(\xi, \phi_i))$$

# Call planner!

$$\theta^+ = \theta - \eta \left[ \underbrace{\nabla_\theta C_\theta(\xi_i^h, \phi_i)}_{\text{(Push down human cost)}} - \underbrace{\nabla_\theta C_\theta(\xi_i, \phi_i)}_{\text{(Push up planner cost)}} \right]$$

# Update cost

# Maximum Entropy Inverse Optimal Control



for  $i = 1, \dots, N$

# Loop over datapoints

$$\xi_i \sim \frac{1}{Z} \exp(-C_\theta(\xi, \phi_i))$$

# Call planner!

$$\theta^+ = \theta - \eta \left[ \nabla_\theta C_\theta(\xi_i^h, \phi_i) - \nabla_\theta C_\theta(\xi_i, \phi_i) \right]$$

(Push down human cost) (Push up planner cost)

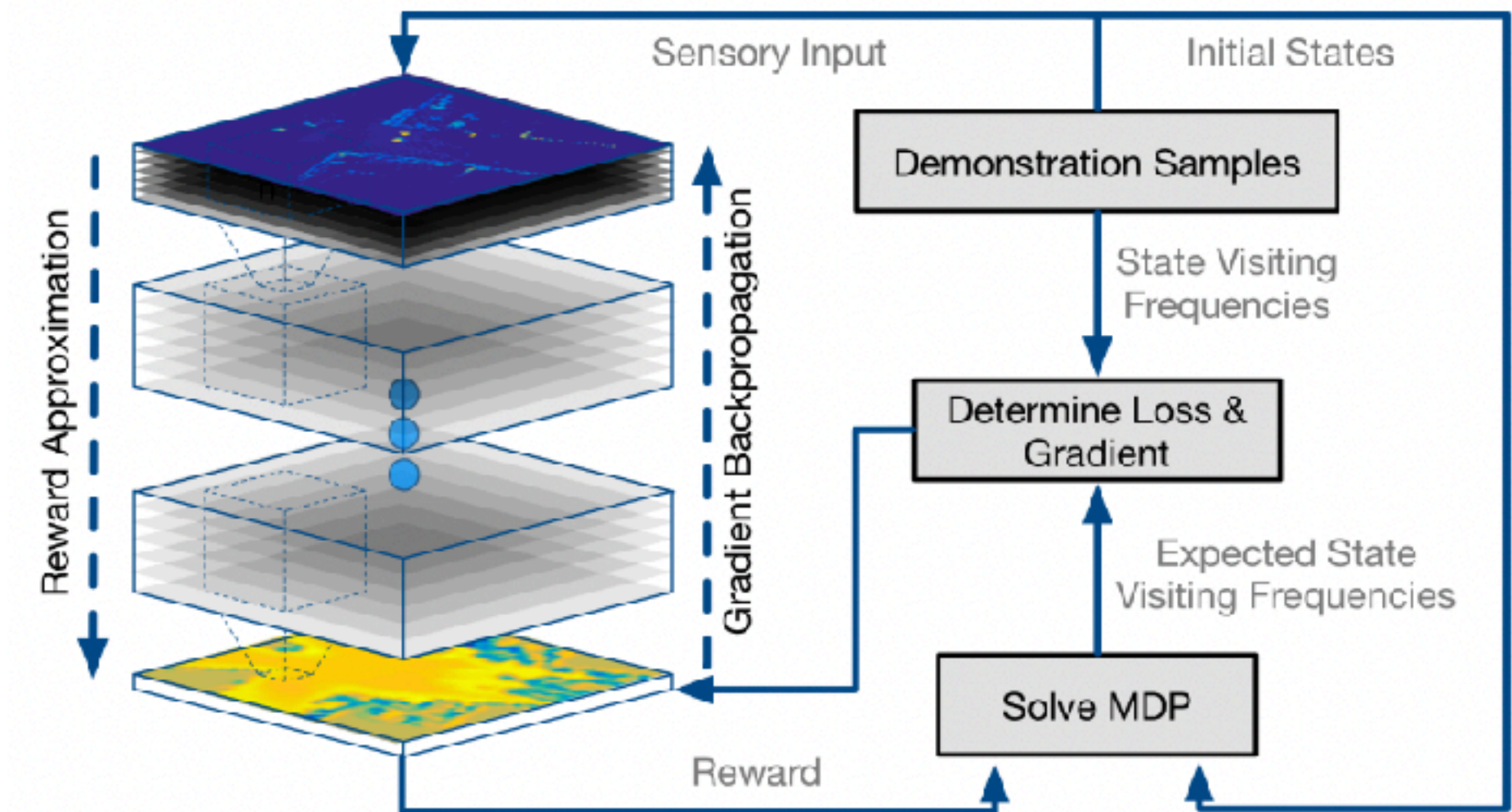
# Update cost

# Deep Max Ent



## Watch This: Scalable Cost-Function Learning for Path Planning in Urban Environments

Markus Wulfmeier<sup>1</sup>, Dominic Zeng Wang<sup>1</sup> and Ingmar Posner<sup>1</sup>

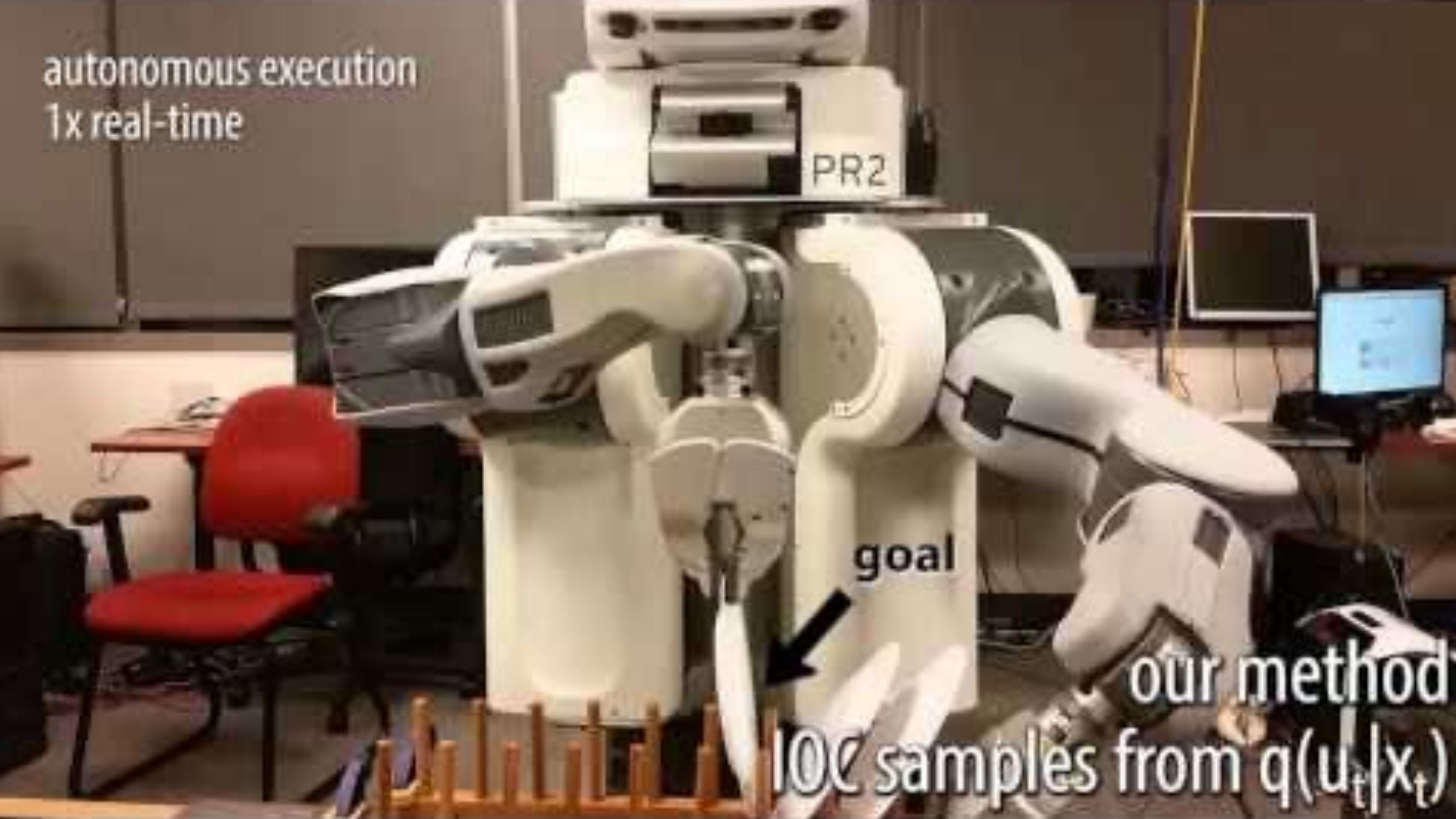


autonomous execution  
1x real-time

PR2

goal

our method  
100 samples from  $q(u_t|x_t)$



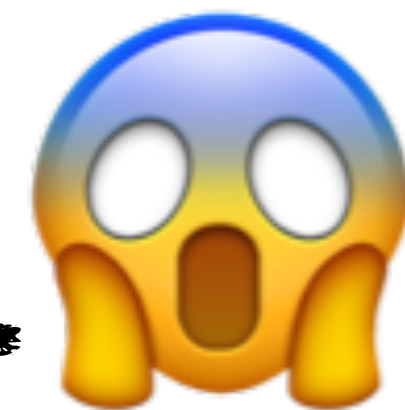
Easy



Medium



Hard



Expert is **realizable**

$$\pi^E \in \Pi$$

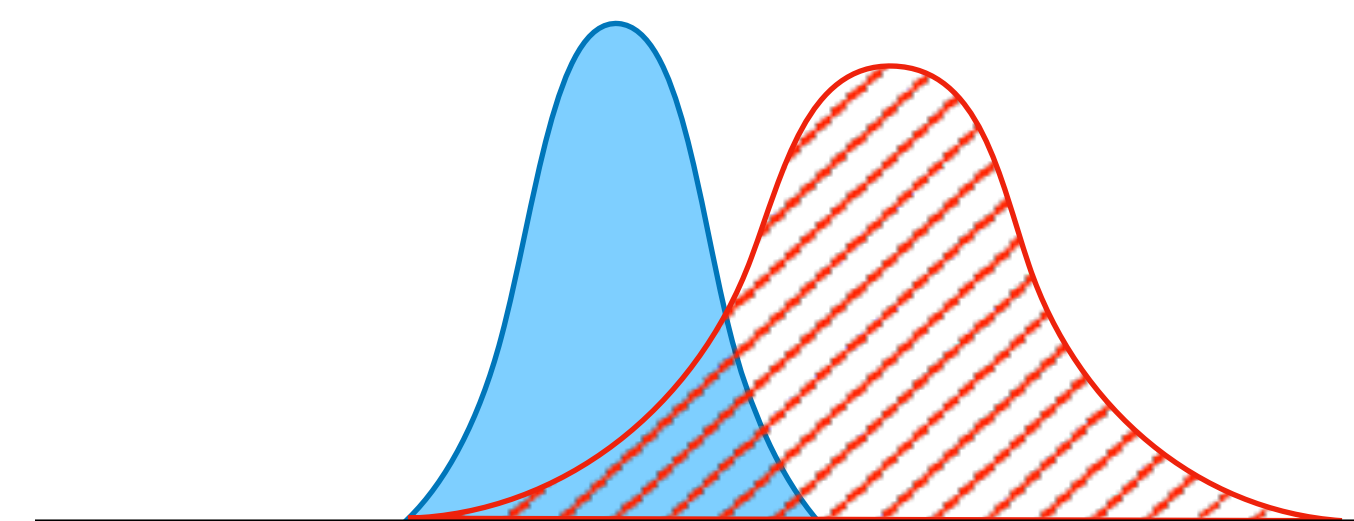
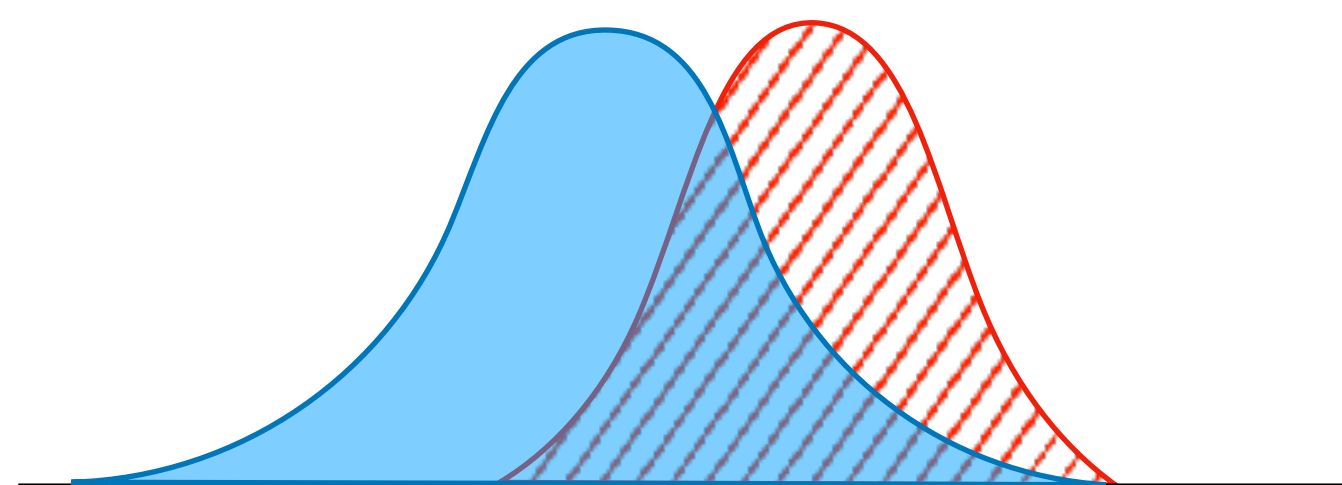
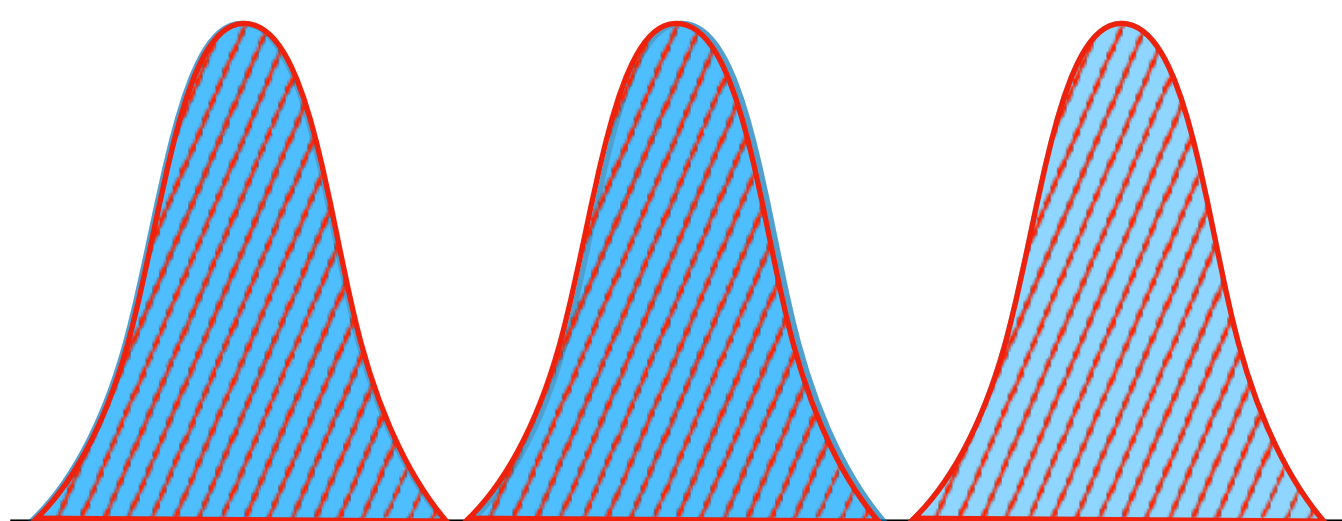
Non-realizable expert  
but full expert support

Non-realizable expert +  
limited expert support

As  $N \rightarrow \infty$ , drive down  
 $\epsilon = 0$  (or Bayes error)

Even as  $N \rightarrow \infty$ ,  
behavior cloning  $O(\epsilon CT)$   
where  $C$  is conc. coeff

Even as  $N \rightarrow \infty$ ,  
behavior cloning  $O(\epsilon T^2)$



Nothing special.

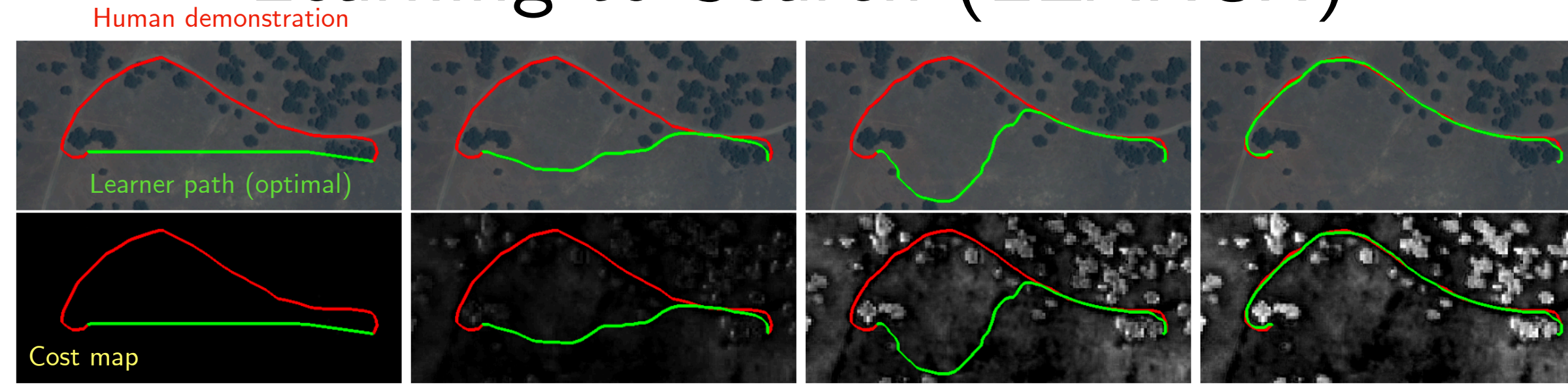
Collect lots of data and  
do Behavior Cloning

Requires **interactive** simulator  
(MaxEntIRL) to match  
distribution  $\Rightarrow O(\epsilon T)$

Requires **interactive** expert  
(DAGGER / **EIL**) to  
provide labels  $\Rightarrow O(\epsilon T)$

# tl;dr

## Learning to Search (LEARCH)



for  $i = 1, \dots, N$


# Loop over datapoints

$$\xi_i^* = \min_{\xi} [C_{\theta}(\xi, \phi_i) - \gamma(\xi, \xi^h)]$$

# Call planner!

$$\theta^+ = \theta - \eta [ \underbrace{\nabla_{\theta} C_{\theta}(\xi_i^h, \phi_i)}_{\text{(Push down human cost)}} - \underbrace{\nabla_{\theta} C_{\theta}(\xi_i^*, \phi_i)}_{\text{(Push up planner cost)}} + \nabla_{\theta} R(\theta) ]$$

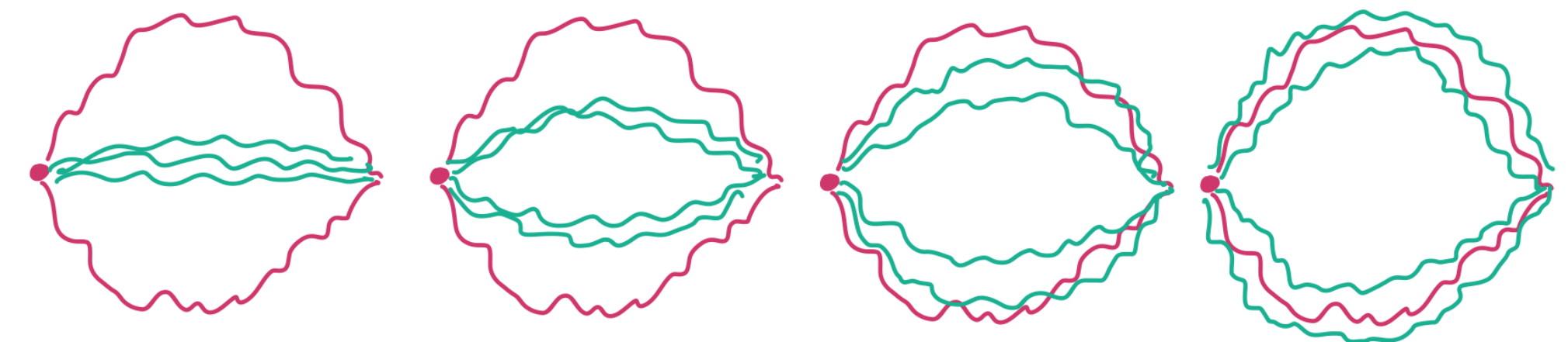
# Update cost



When the expert is  
Suboptimal  
Noisy  
Privileged Information

LEARCH does NOT converge!!

## Maximum Entropy Inverse Optimal Control



for  $i = 1, \dots, N$

# Loop over datapoints

$$\xi_i \sim \frac{1}{Z} \exp(-C_{\theta}(\xi, \phi_i))$$

# Call planner!

$$\theta^+ = \theta - \eta [ \underbrace{\nabla_{\theta} C_{\theta}(\xi_i^h, \phi_i)}_{\text{(Push down human cost)}} - \underbrace{\nabla_{\theta} C_{\theta}(\xi_i, \phi_i)}_{\text{(Push up planner cost)}} ]$$

# Update cost