

Imitation Learning by Inverse Reinforcement Learning

Jan Peters
Gerhard Neumann

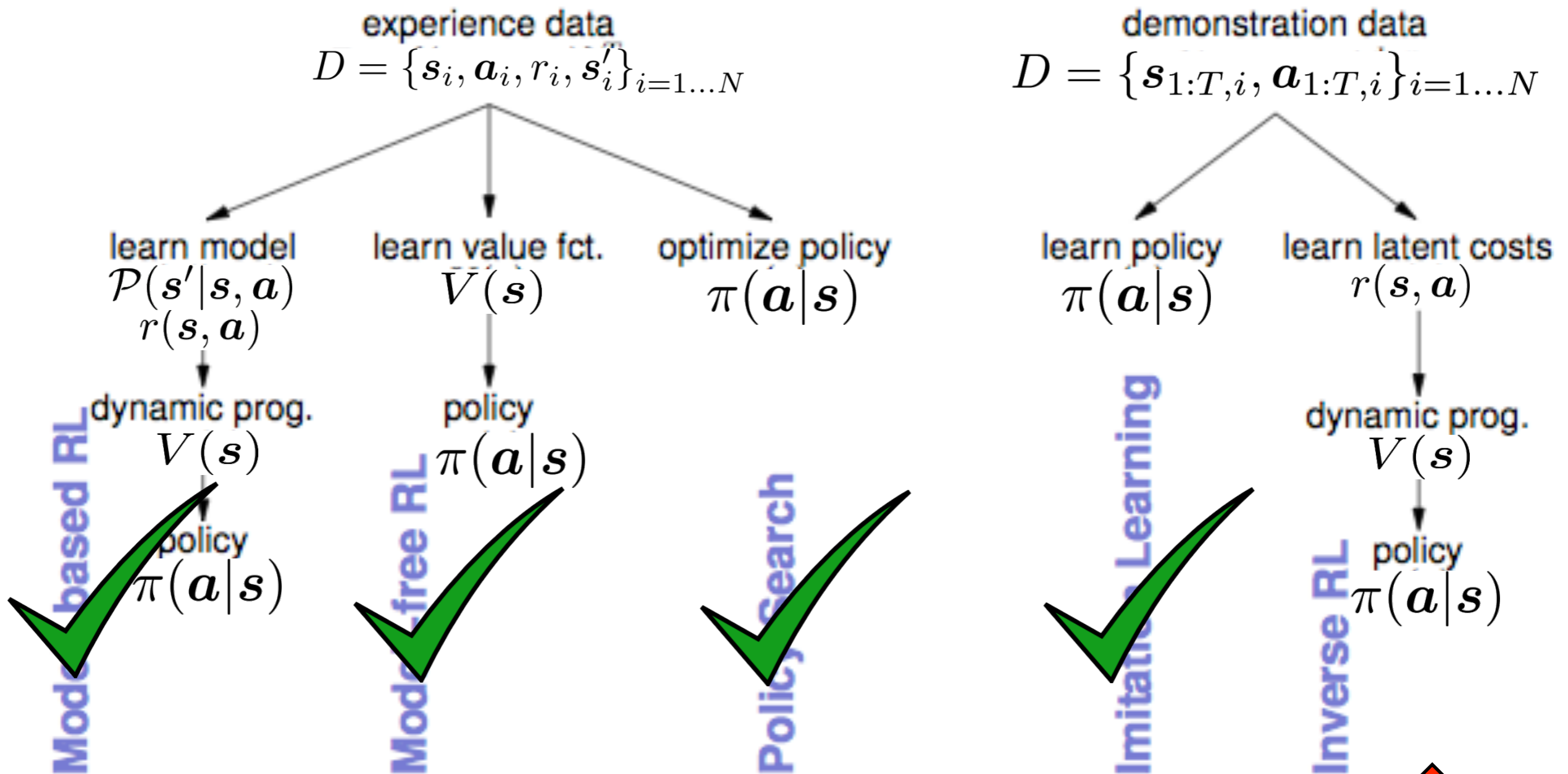
Inspired by Slides from P. Abbeel,
Drew Bagnell and others



Purpose of this Lecture

- Learn an alternative approach to imitation learning
- What is the best way to imitate a teacher?
 - Learn its policy? Behavioral cloning
 - Needs a lot of demonstrations to generalize the behavior
 - Learn its intention / goals? Inverse Reinforcement Learning
 - Inverse Optimal Control, Inverse Optimal Planning
 - Determine the cost function of the teacher in order to obtain optimal behavior.
 - More concise description of behavior

Bigger Picture





Outline of the Lecture

1. Introduction

- I. Comparison to Behavioral Cloning

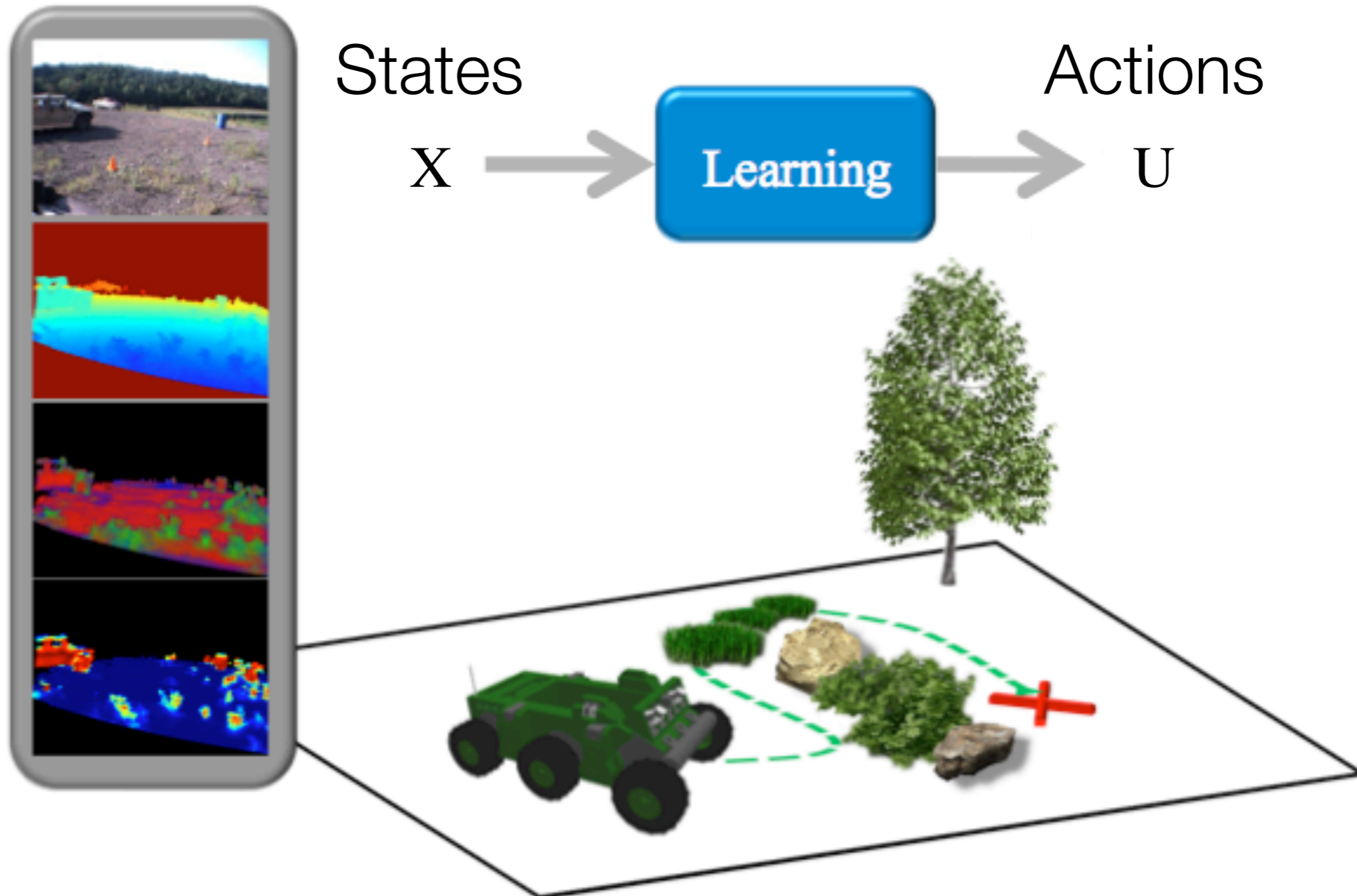
2. Categories of IRL

- I. Maximum Margin
- II. Feature Matching by Max. Entropy
- III. Policy parametrized by rewards

3. Applications

4. Conclusion

Behavioral Cloning



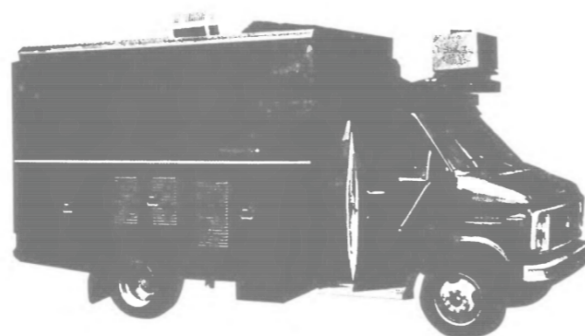


What may be wrong here? Remember ALVINN?

There are difficulties involved with training “on-the-fly” with real images. If the network is not presented with sufficient variability in its training exemplars to cover the conditions it is likely to encounter when it takes over driving from the human operator, it will not develop a sufficiently robust representation and will perform poorly. In addition, the network must not solely be shown examples of accurate driving, but also how to recover (i.e. return to the road center) once a mistake has been made. Partial initial training on

Disadvantages of Direct Imitation Learning

- Needs a lot of demonstrations to generalize
- High variability in the demonstrations
- Demonstrate how to recover from mistakes





Motivation for inverse RL

Apprenticeship learning/Imitation learning through inverse RL

Presupposition: reward function provides the most succinct and transferable definition of the task

Has enabled advancing the state of the art in various robotic domains

Modeling of other agents, both adversarial and cooperative

Scientific questions

Model animal and human behavior

E.g., bee foraging, songbird vocalization. [See intro of Ng and Russell, 2000 for a brief overview.]

Meet Crusher...



More Crusher pictures...

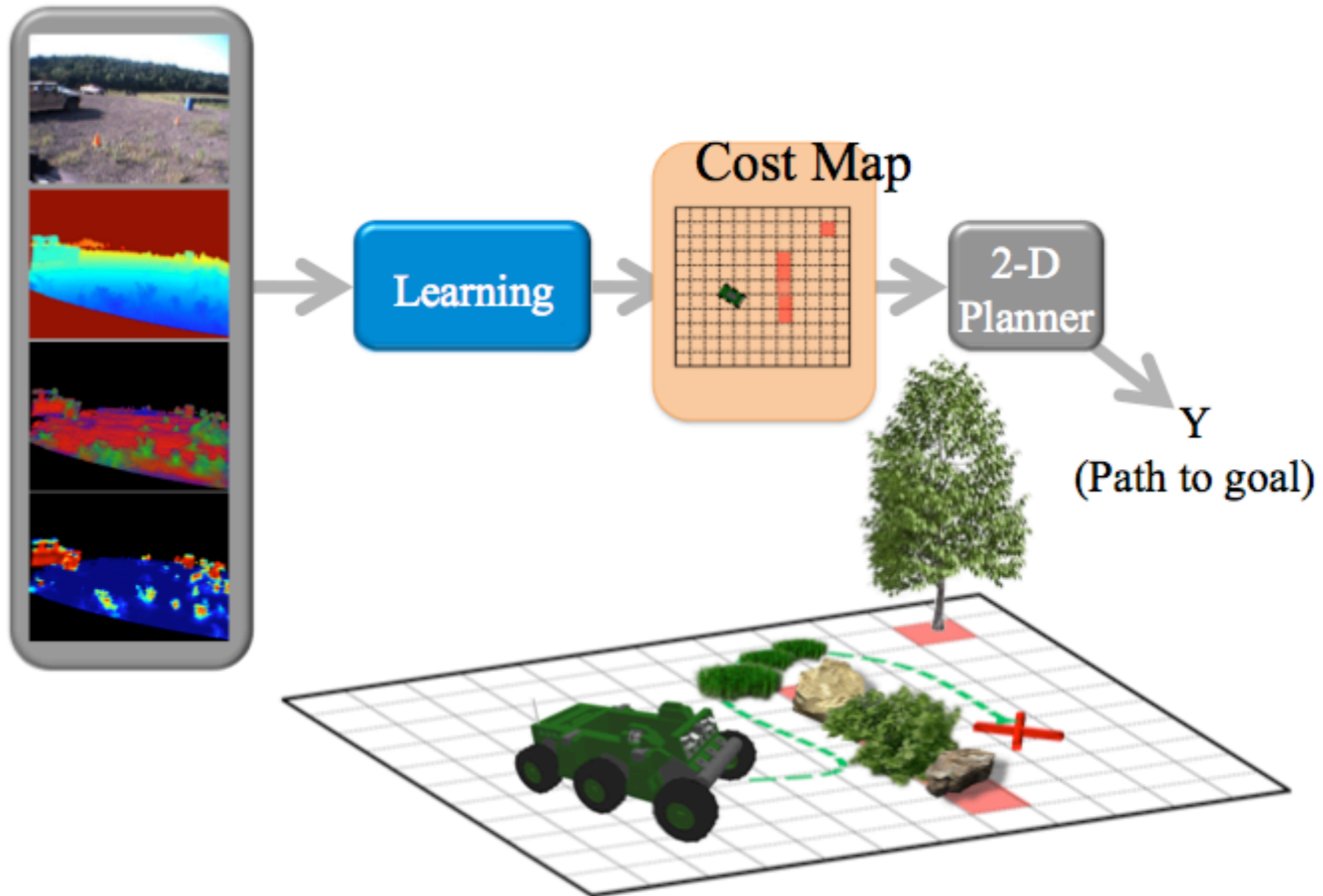


More Crusher pictures...

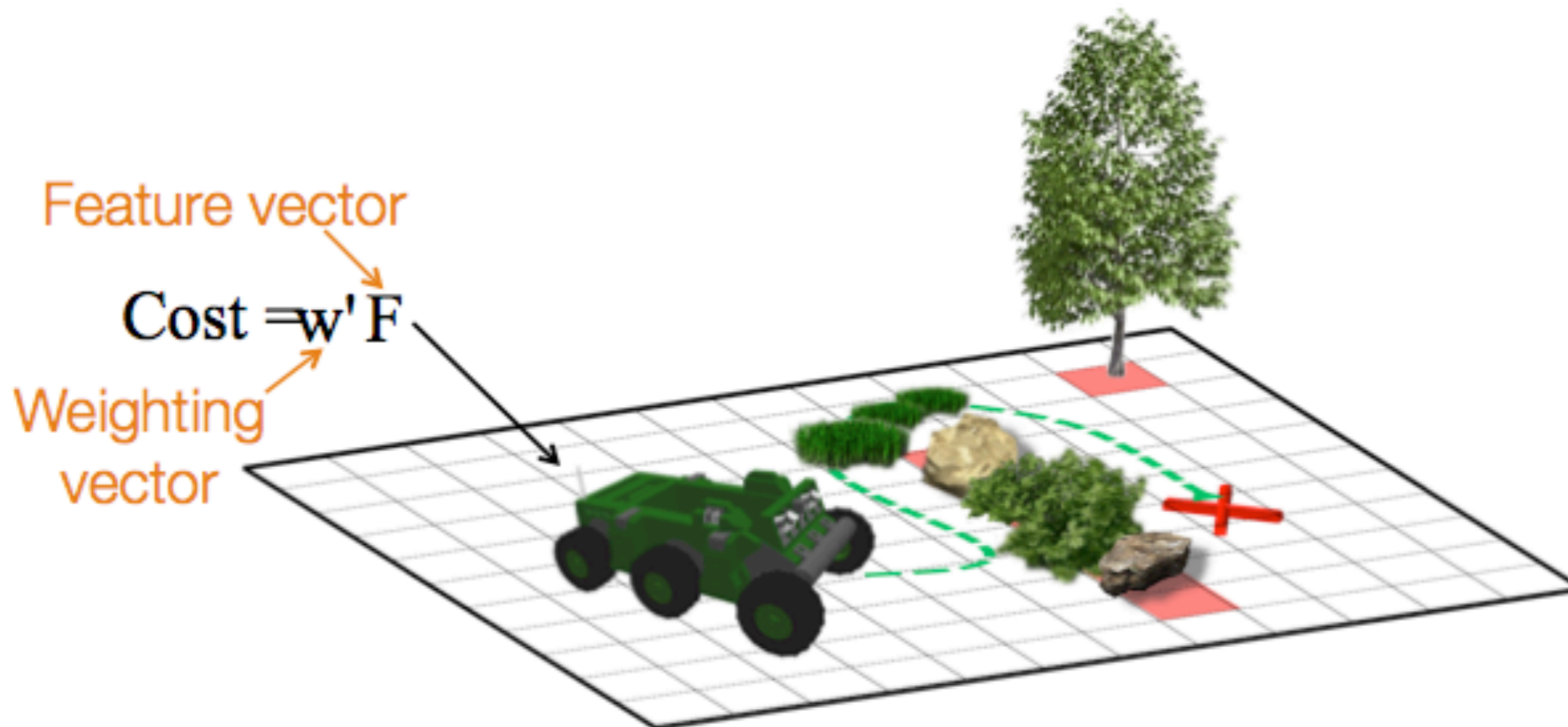




Inverse Reinforcement Learning

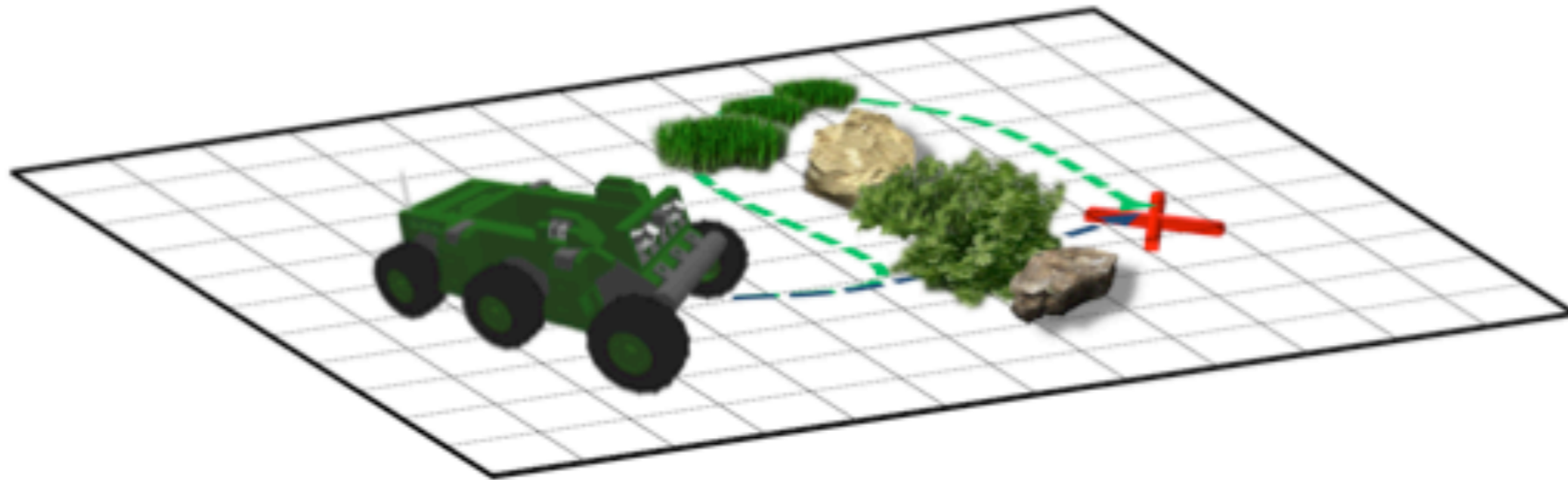
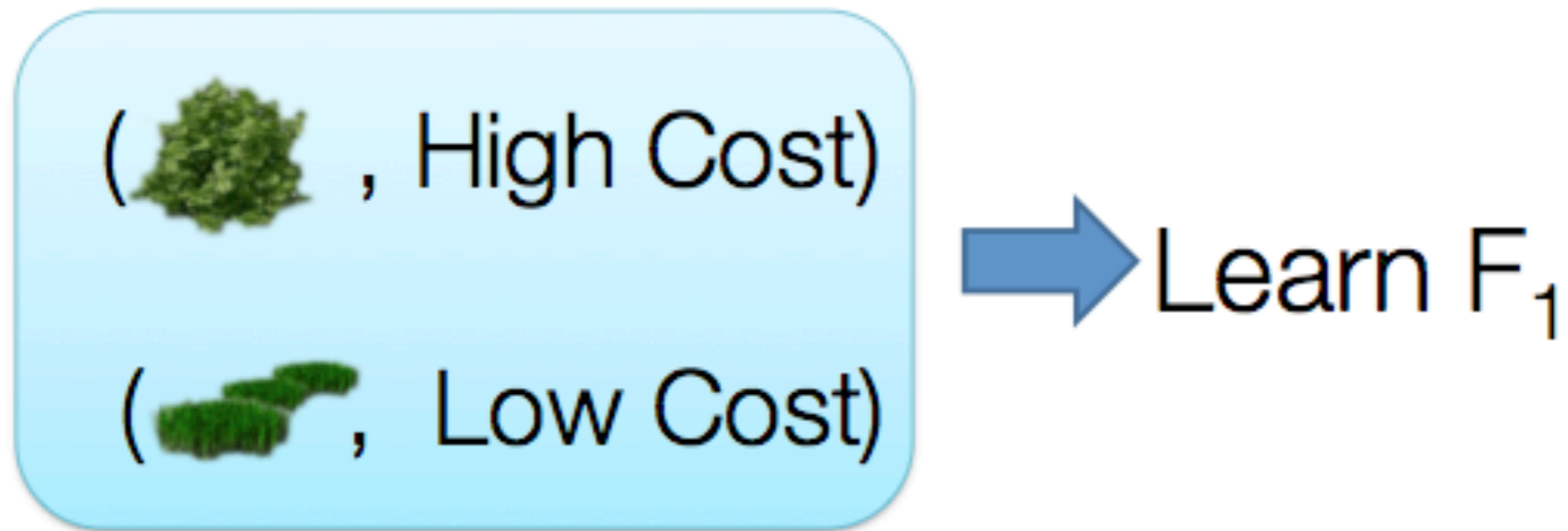


Recovering Cost!

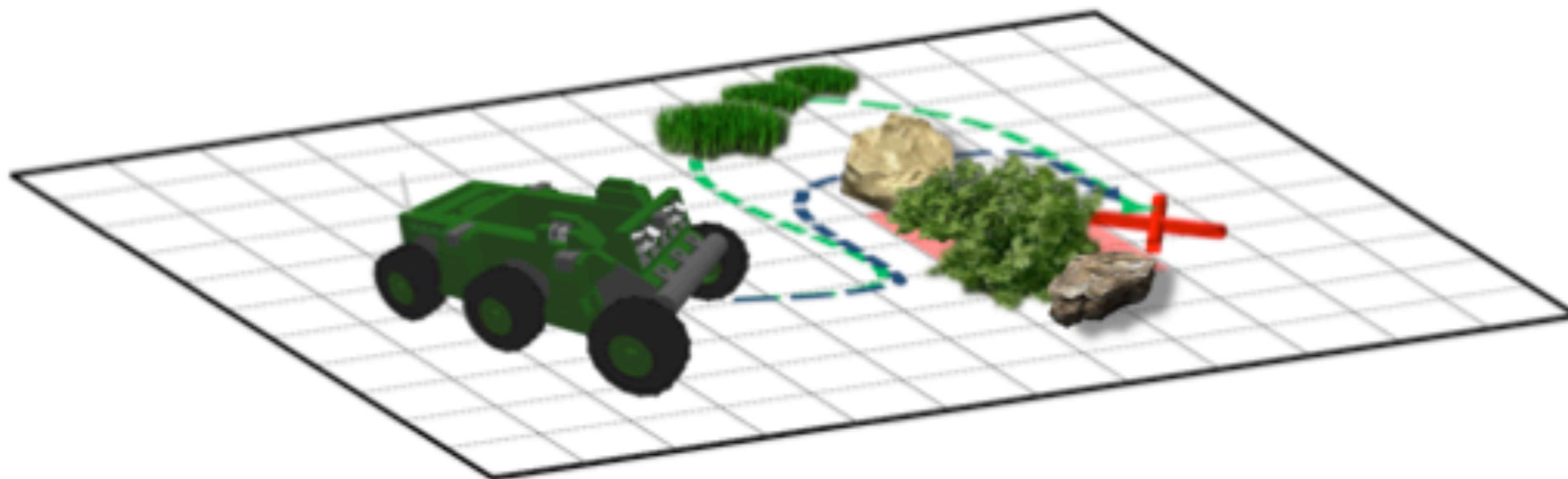
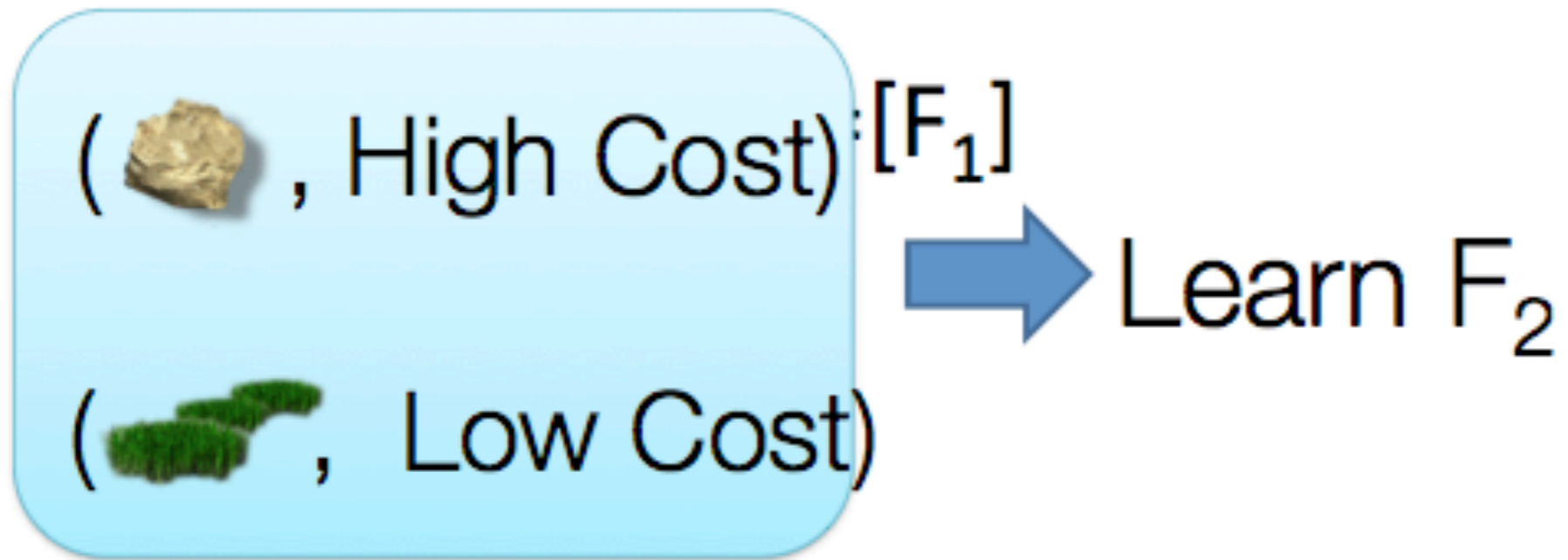
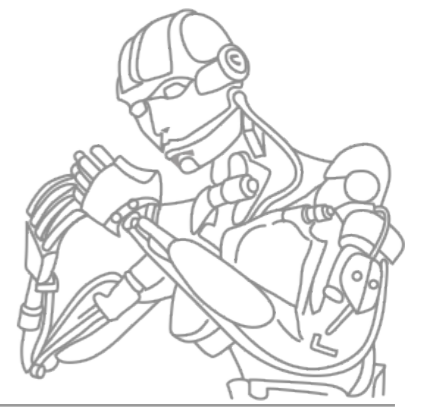


Ratliff, Bagnell, Zinkevich 2005
Ratliff, Bradley, Bagnell, Chestnutt, NIPS 2006
Silver, Bagnell, Stentz, RSS 2008

Recovering Cost!



Recovering Cost!



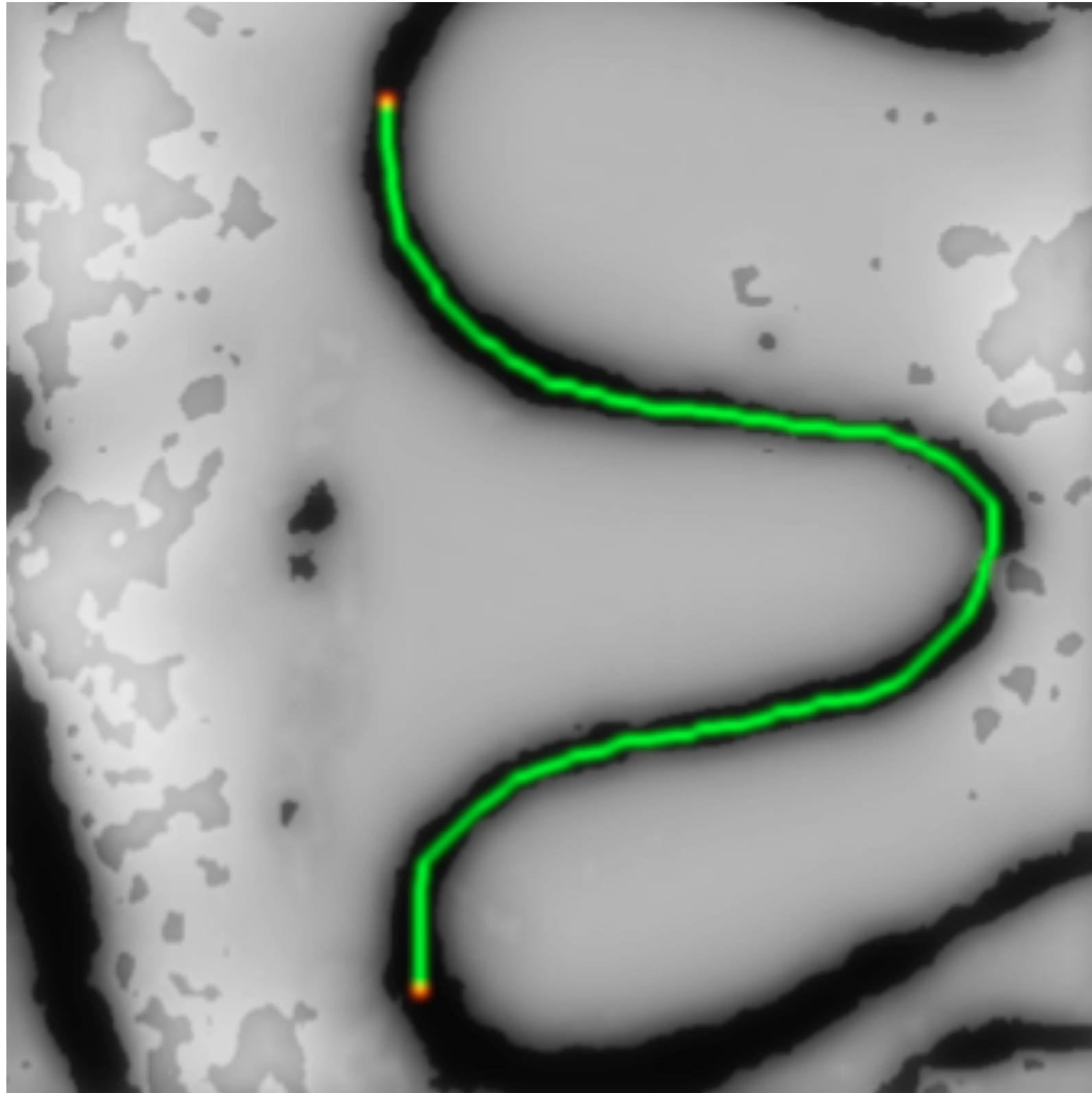
Collect paths by teleoperation



Training: Stay on the road



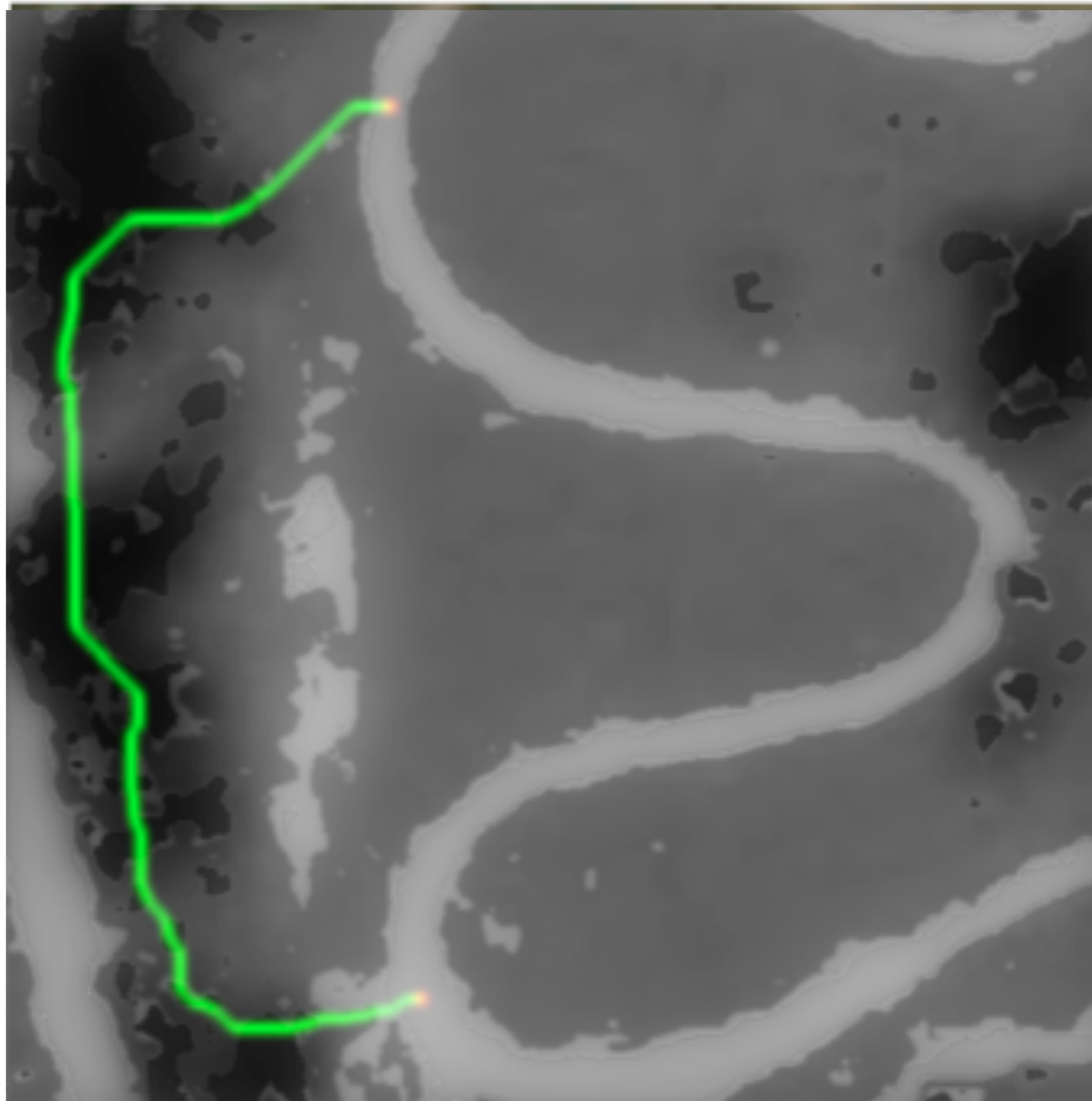
Test: Stay on the road



Training: Avoid the Road



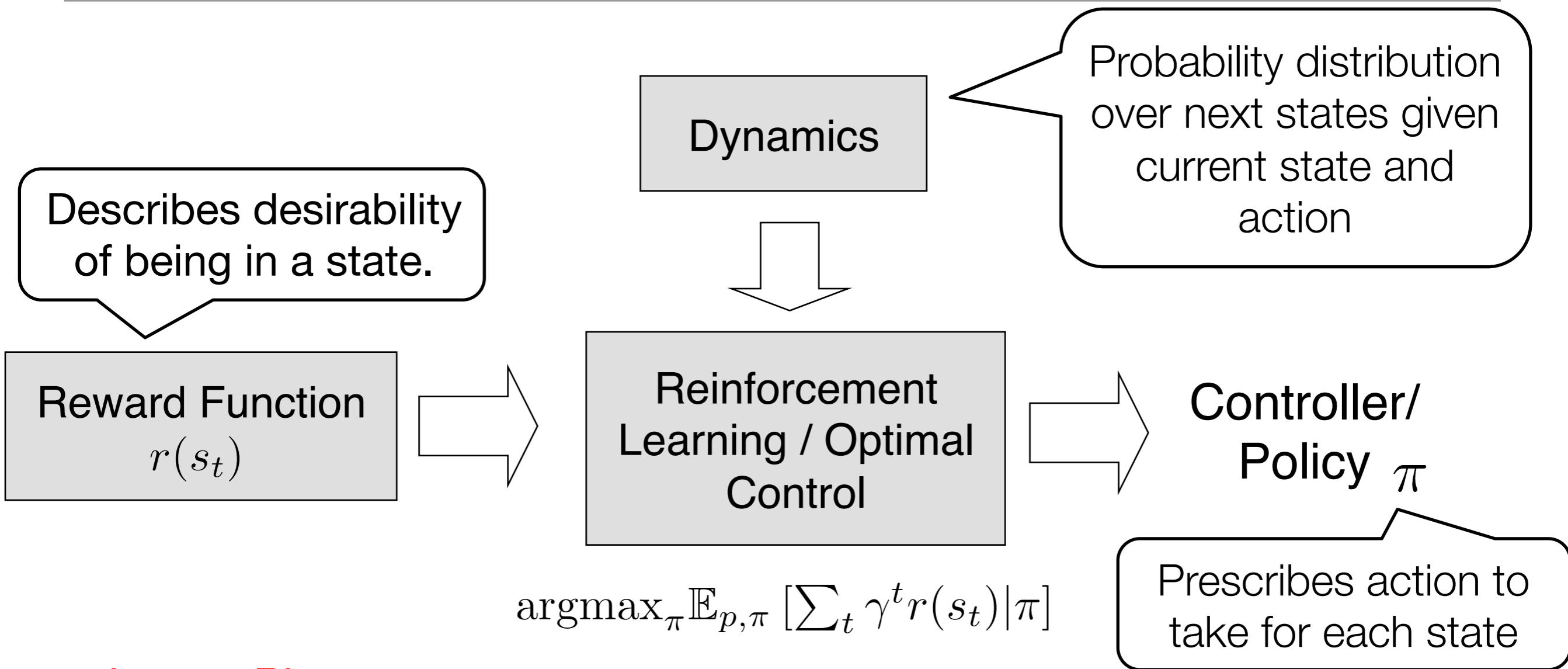
Test: Avoid the Road





example path

High-level picture



Inverse RL:

- Given **Policy** and **Model**, can we recover **R**?
- More generally, given execution traces, can we recover **r**?



Outline of the Lecture

1. Introduction

- I. Comparison to Behavioral Cloning

2. Categories of IRL

- I. Maximum Margin
- II. Feature Matching by Max. Entropy
- III. Policy parametrized by rewards

3. Applications

4. Conclusion



Problem setup: Behavioral Cloning

Input:

Teacher's demonstrations: $D = \{\mathbf{s}_{1:T,i}, \mathbf{a}_{1:T,i}\}_{i=1\dots N}$

Trace of the teacher's policy $\pi^*(\mathbf{a}|\mathbf{s})$

And its "long-term behavior" $\mu^*(\mathbf{s})$

Formulated as standard machine learning problem

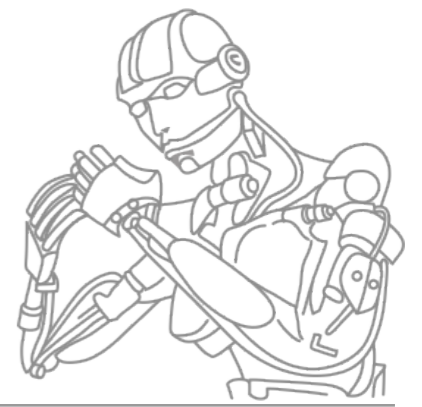
Fix a policy class (neural network, decision tree, deep belief net, dynamical systems, ...)

Estimate a policy $\hat{\pi}(\mathbf{a}|\mathbf{s})$ from the training examples D

Problem:

There will always be an error in the estimation of the policy

Small error in the policy \Rightarrow possibly large error in long-term behavior $\hat{\mu}(\mathbf{s})$



Problem setup: Inverse RL

Input:

Teacher's demonstrations: $D = \{\mathbf{s}_{1:T,i}, \mathbf{a}_{1:T,i}\}_{i=1\dots N}$

Trace of the teacher's policy $\pi^*(\mathbf{a}|\mathbf{s})$

And its "long-term behavior" $\mu^*(\mathbf{s})$

State and Action Space

Transition model: $p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$

No reward function $r(\mathbf{s}_t)$

Inverse RL:

Can we recover $r(\mathbf{s}_t)$ that explains the policy $\pi^*(\mathbf{a}|\mathbf{s})$ (and its long-term behavior $\mu^*(\mathbf{s})$?)

Apprenticeship Learning

Can we use $r(\mathbf{s}_t)$ to obtain a policy $\hat{\pi}(\mathbf{a}|\mathbf{s})$?

Inverse RL vs. Behavioral Cloning



Behavioral Cloning:

Simple to implement

No assumptions on the model/MDP

We might not reproduce the long term behavior

Representation: Policy

Hard to generalize

Needs many samples

Inverse RL:

Requires Planning / Solving an MDP

Hard for many interesting MDPs (e.g. high-DoF robots)

Representation: Reward

Compact description

Easy to transfer to new tasks

Basic principle



Find a reward function $r^*(s_t)$ which explains the expert behavior

Assume expert is optimal w.r.t. to $r^*(s_t)$

I.e., find $r^*(s_t)$ such that

$$\mathbb{E}_{p,\pi} [\sum_t \gamma^t r^*(s_t) | \pi^*] \geq \mathbb{E}_{p,\pi} [\sum_t \gamma^t r^*(s_t) | \pi], \forall \pi$$

In fact a convex feasibility problem, but many challenges:

1. **Ill-posed:** $r^*(s_t) = 0$ is a solution, reward function ambiguity
2. **Limited Data:** We typically only observe expert traces rather than the entire expert optimal policy --- how to compute left-hand side?
3. **Optimality Assumption:** Assumes the expert is indeed optimal --- otherwise infeasible

- 27 4. **Computation:** assumes we can enumerate all policies



Outline of the Lecture

1. Introduction

- I. Comparison to Behavioral Cloning

2. Categories of IRL

- I. Maximum Margin
- II. Feature Matching by Max. Entropy
- III. Policy parametrized by rewards

3. Applications

4. Conclusion

Feature based reward function



Lets assume the reward function is linear in some features,

I.e. $r(\mathbf{s}) = \mathbf{w}^T \phi(\mathbf{s})$, where $\phi(\mathbf{s})$ is a n -dimensional feature vector

$$\begin{aligned}\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t) | \pi \right] &= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{w}^T \phi(\mathbf{s}_t) | \pi \right] \\ &= \mathbf{w}^T \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \phi(\mathbf{s}_t) | \pi \right] \\ &= \mathbf{w}^T \psi(\pi)\end{aligned}$$

where $\psi(\pi)$ is the expected discounted feature vector of policy π

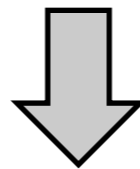
Subbing into: $\mathbb{E}_{p, \pi^*} [\sum_t \gamma^t r^*(s_t) | \pi^*] \geq \mathbb{E}_{p, \pi} [\sum_t \gamma^t r^*(s_t) | \pi], \forall \pi$

gives us: Find \mathbf{w}^* such that $\mathbf{w}^{*T} \psi(\pi^*) \geq \mathbf{w}^{*T} \psi(\pi), \forall \pi$

Feature based reward function



$$\mathbb{E}_{p,\pi} [\sum_t \gamma^t r^*(s_t) | \pi^*] \geq \mathbb{E}_{p,\pi} [\sum_t \gamma^t r^*(s_t) | \pi], \forall \pi$$



Find \mathbf{w}^* such that $\mathbf{w}^{*T} \psi(\pi^*) \geq \mathbf{w}^{*T} \psi(\pi), \forall \pi$

Feature expectations can be readily estimated from sample trajectories

Solves limited data challenge

The number of expert demonstrations required scales with the number of features in the reward function.

The number of expert demonstration required does not depend on

Complexity of the expert's optimal policy

Size of the state space

Basic principle



Find a reward function $r^*(s_t)$ which explains the expert behavior

Assume expert is optimal w.r.t. to $r^*(s_t)$

I.e., find $r^*(s_t)$ such that

$$\mathbb{E}_{p,\pi} [\sum_t \gamma^t r^*(s_t) | \pi^*] \geq \mathbb{E}_{p,\pi} [\sum_t \gamma^t r^*(s_t) | \pi], \forall \pi$$

In fact a convex feasibility problem, but many challenges:

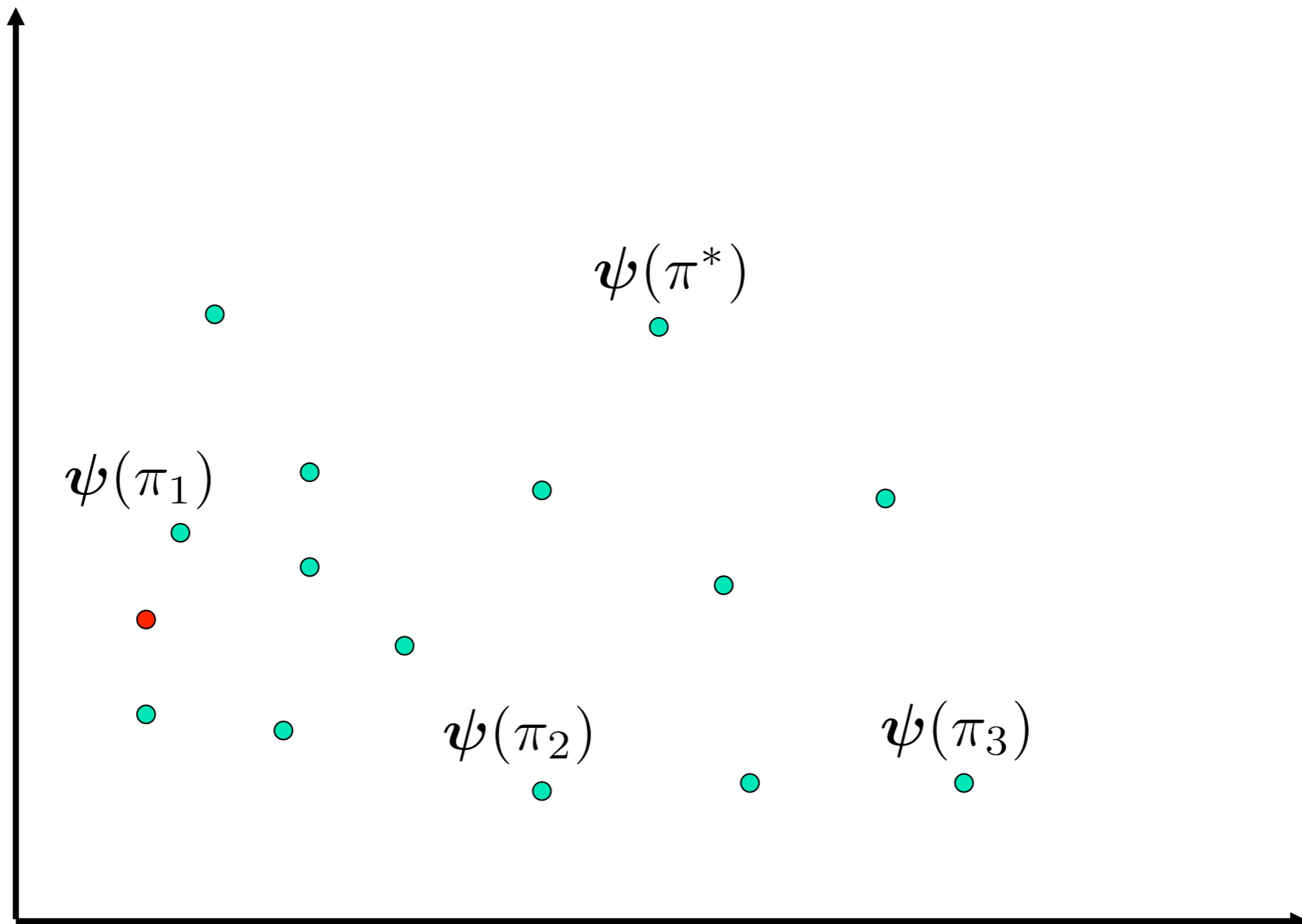
1. **Ill-posed:** $r^*(s_t) = 0$ is a solution, reward function ambiguity
2. **Limited Data:** We typically only observe expert traces rather than the entire expert optimal policy --- how to compute left-hand side?
3. **Optimality Assumption:** Assumes the expert is indeed optimal --- otherwise infeasible

- 31 4. **Computation:** assumes we can enumerate all policies



Constraint generation

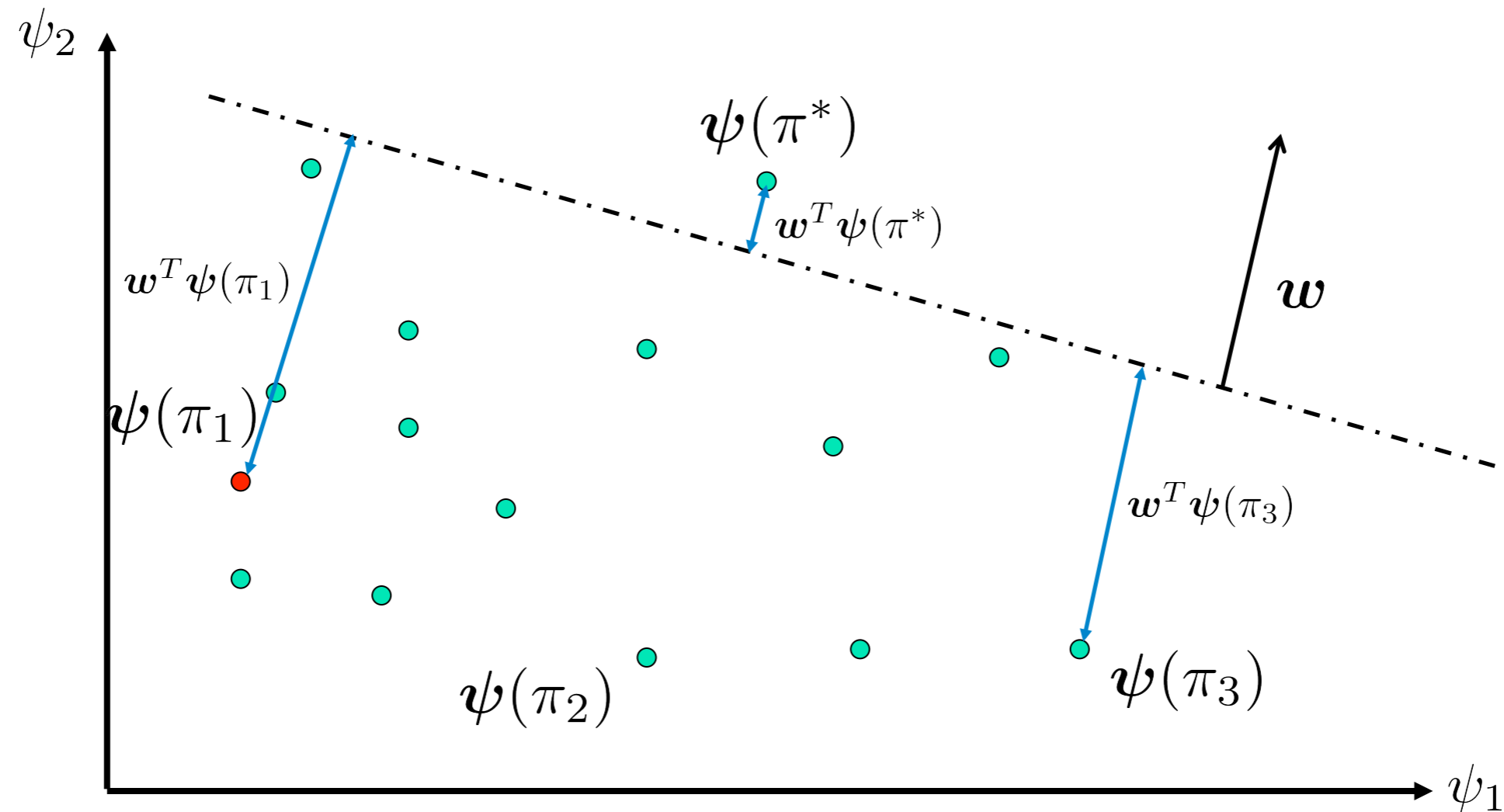
Every policy has a corresponding feature expectation vector, which for visualization purposes we assume to be 2D



Constraint generation



Linear parametrization: We need to find a separating hyper-plane given by w



Scalar product $w^T \psi$ gives us the (positive or negative) distance to the separating hyper-plane



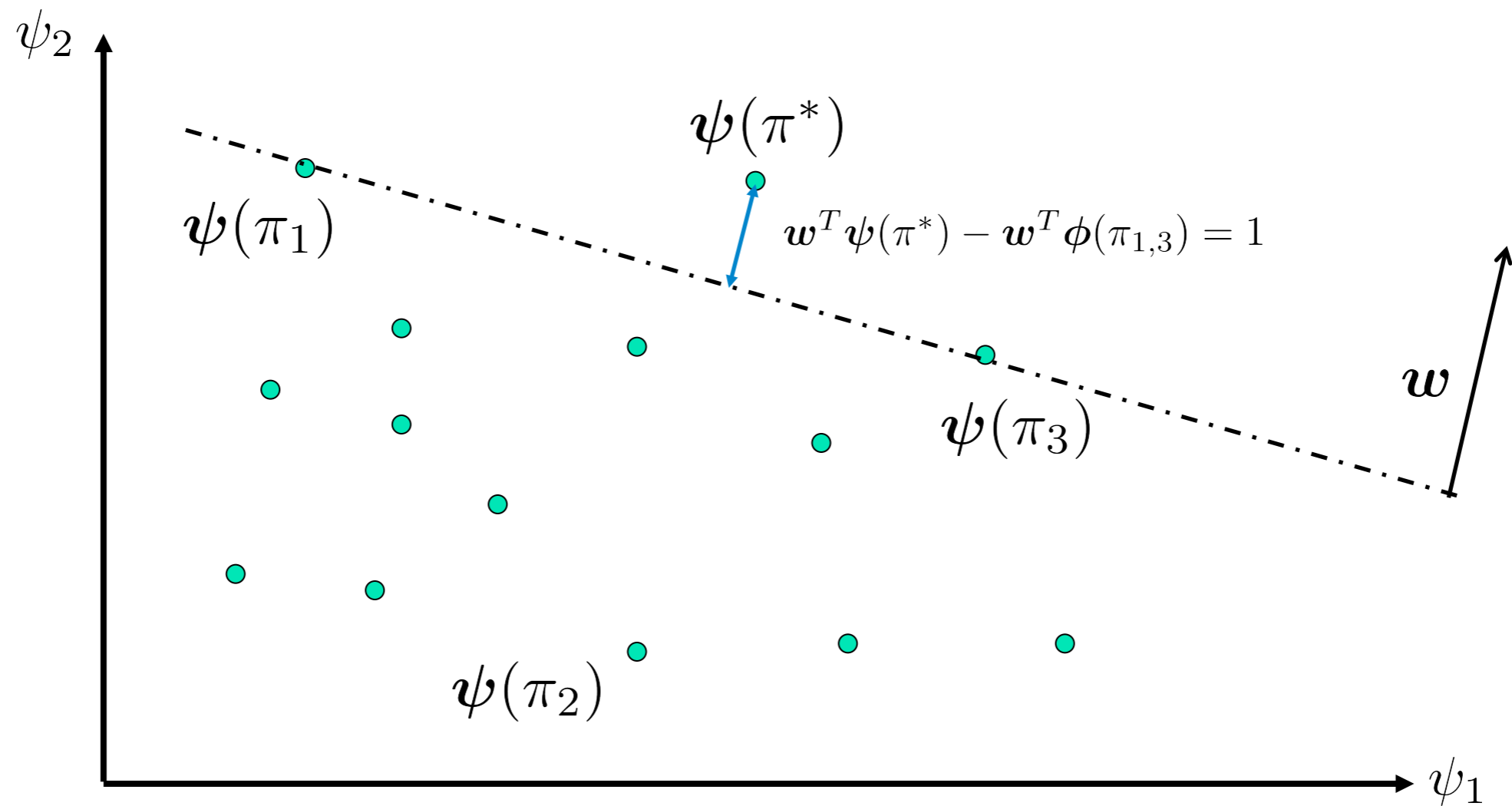
Ill-posed Problem

Standard max margin:

Smallest weight vector with predefined reward margin of 1

$$\begin{array}{ll} \min_{\mathbf{w}} & \mathbf{w}^T \mathbf{w} \\ \text{s.t.} & \mathbf{w}^T \boldsymbol{\psi}(\pi^*) \geq \mathbf{w}^T \boldsymbol{\psi}(\pi) + 1, \quad \forall \pi \end{array}$$

Max. margin solution



Similar interpretation as a support vector machine



Ill-posed Problem

Structured max margin:

Smallest weight vector

Margin depends on difference of policies $m(\pi^*, \pi)$

$$\begin{aligned} \min_{\mathbf{w}} \quad & \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & \mathbf{w}^T \psi(\pi^*) \geq \mathbf{w}^T \psi(\pi) + m(\pi^*, \pi), \quad \forall \pi \end{aligned}$$

Justification: margin should be larger for policies that are very different from π^* .

Example for $m(\pi^*, \pi)$:

Sum of minimum distances from generated path the example path

Basic principle



Find a reward function $r^*(s_t)$ which explains the expert behavior

Assume expert is optimal w.r.t. to $r^*(s_t)$

I.e., find $r^*(s_t)$ such that

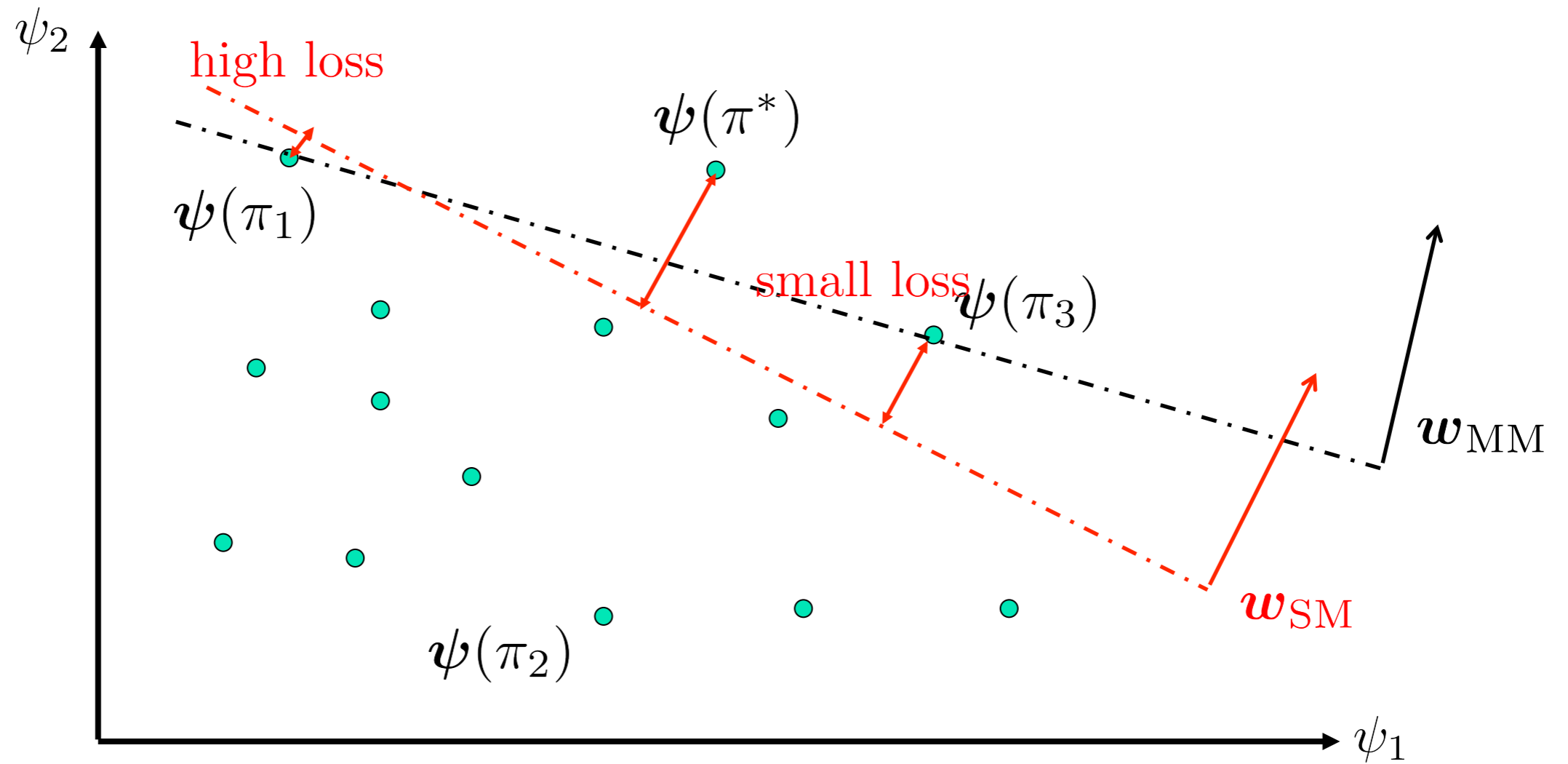
$$\mathbb{E}_{p,\pi} [\sum_t \gamma^t r^*(s_t) | \pi^*] \geq \mathbb{E}_{p,\pi} [\sum_t \gamma^t r^*(s_t) | \pi], \forall \pi$$

In fact a convex feasibility problem, but many challenges:

1. Ill-posed: $r^*(s_t) = 0$ is a solution, reward function ambiguity
2. Limited Data: We typically only observe expert traces rather than the entire expert optimal policy --- how to compute left-hand side?
3. **Optimality Assumption:** Assumes the expert is indeed optimal --- otherwise infeasible

- 37 4. **Computation:** assumes we can enumerate all policies

Structured max margin solution





Expert suboptimality

Structured prediction max margin with slack variables:

Every constraint can be violated a bit

Minimize the amount of violation

$$\min_{\mathbf{w}, \xi} \quad \mathbf{w}^T \mathbf{w} + C\xi$$

$$\text{s.t.} \quad \mathbf{w}^T \boldsymbol{\psi}(\pi^*) \geq \mathbf{w}^T \boldsymbol{\psi}(\pi) + m(\pi^*, \pi) - \xi, \quad \forall \pi$$

Easy to extend to multiple MDPs

Resolved: access to π^* , ambiguity, expert suboptimality

One challenge remains: very large number of constraints

Ratliff et. al. use subgradient methods.

In this lecture: constraint generation

Constraint generation



Initialize $\Pi = \{\}$ **and then iterate** $k = 1 \dots$

Solve

$$\mathbf{w}^{(k)} = \underset{\mathbf{w}}{\operatorname{argmin}} \min_{\xi} \mathbf{w}^T \mathbf{w} + C\xi$$

$$\text{s.t.} \quad \mathbf{w}^T \boldsymbol{\psi}(\pi^*) \geq \mathbf{w}^T \boldsymbol{\psi}(\pi^{(i)}) + m(\pi^*, \pi^{(i)}) - \xi, \quad \forall \pi^{(i)} \in \Pi$$

Find the most violated constraint

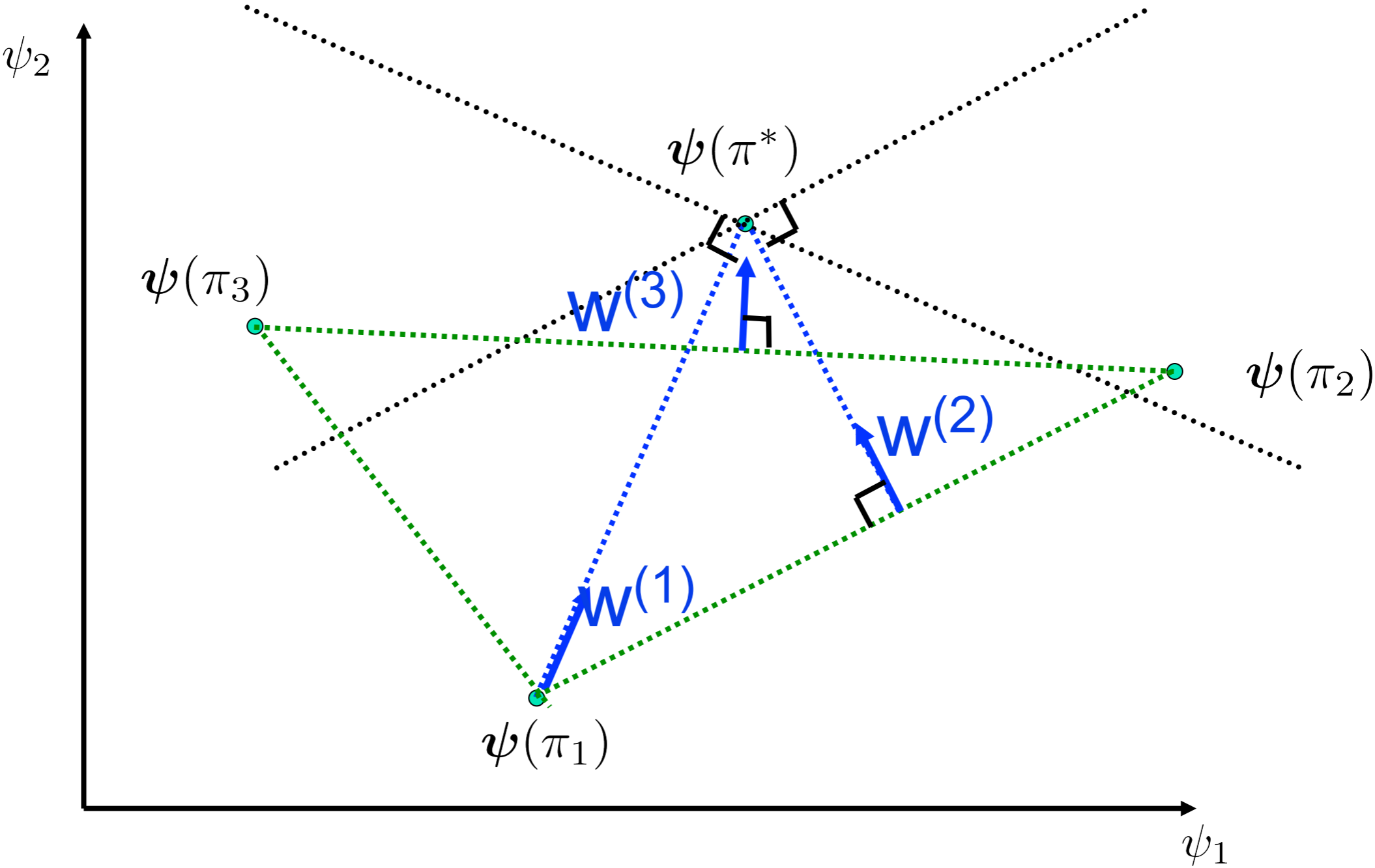
$$\pi^{(k)} = \max_{\pi} \mathbf{w}^{(k),T} \boldsymbol{\psi}(\pi) + m(\pi^*, \pi)$$

Compute optimal policy for the current estimate of the reward function (+ loss augmentation m), e.g., dynamic programming

Add $\pi^{(k)}$ **to set of policies** Π

If no constraint violations were found, we are done.

Algorithm example run



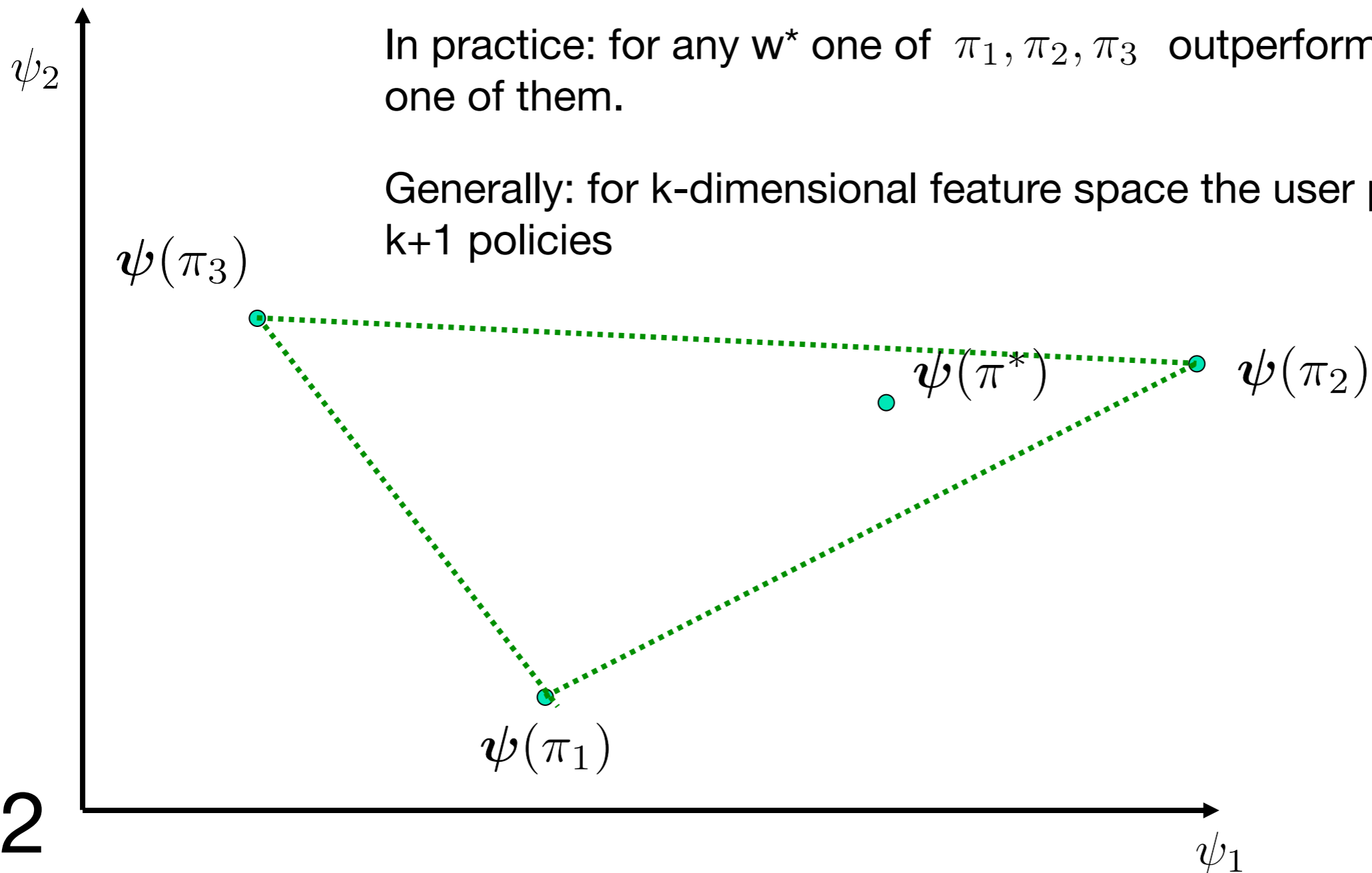
Suboptimal expert case



Can match expert by stochastically mixing between 3 policies

In practice: for any w^* one of π_1, π_2, π_3 outperforms $\pi^* \rightarrow$ pick one of them.

Generally: for k -dimensional feature space the user picks between $k+1$ policies





Outline of the Lecture

1. Introduction

- I. Comparison to Behavioral Cloning

2. Categories of IRL

- I. Maximum Margin
- II. Feature Matching by Max. Entropy
- III. Policy parametrized by rewards

3. Applications

4. Conclusion

Classical Approach to Statistical Modeling: "The Principle of Maximum Entropy"



Premise: Statistical modeling should be performed with the least commitment possible.

- Predict using the probability distribution **minimally committed/ maximally uncertain/highest Shannon entropy....**
- ...subject to it agrees with known constraints
- “almost all” distributions are close to the MaxEnt one
- **Uncertainty allows proper treatment of suboptimal demonstrations!**



Maximum-entropy approach to inference



What is the **maximum entropy distribution** with given 1st and 2nd order and moments?

$$\operatorname{argmax}_p \underbrace{- \sum_x p(x) \log p(x)}_{\text{Entropy}}$$

$$\text{s.t.} \quad \sum_x p(x)x = m_1, \quad \sum_x p(x)x^2 = m_2, \quad \sum_x p(x) = 1$$

Solution:

$$p(x) \propto \exp(\lambda_1 x + \lambda_2 x^2), \quad \text{where } \lambda_{1,2} \text{ are lagrangian multipliers}$$

That's a Gaussian! ... which is of course well known that a Gaussian has maximum entropy of all distributions with given mean and variance

Maximum-entropy approach to IRL [Ziebart 2008]



Maximize entropy over paths with a given feature expectation

$$\begin{aligned} \operatorname{argmax}_p \quad & - \sum_{\tau} p(\tau) \log p(\tau) \\ \text{s.t.} \quad & \sum_{\tau} p(\tau) \psi(\tau) = \psi(\pi^*), \quad \sum_{\tau} p(\tau) = 1 \end{aligned}$$

Solution: $p(\tau) \propto \exp(\mathbf{w}^T \psi(\tau))$

w is now a lagrangian multiplier

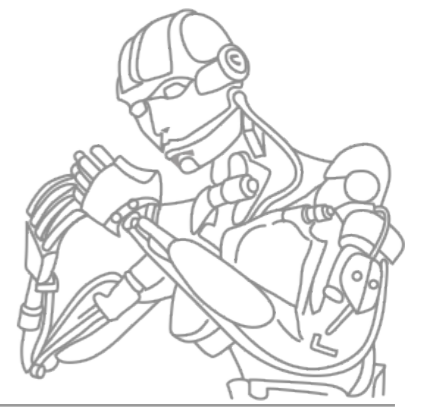
We obtain a soft-max distribution over trajectories

Return of the trajectories: $R(\tau) = \mathbf{w}^T \sum_t \gamma^t \phi(\mathbf{s}_t) = \mathbf{w}^T \psi(\tau)$

Problem: Does not take system dynamics into account

Trajectory could have huge return, but is very unlikely due to system dynamics

Maximum-Causal-entropy IRL [Ziebart 2010]



Maximize entropy of the policy with a given feature expectation

$$\begin{aligned} \operatorname{argmax}_{\pi} \quad & - \sum_{t,a} \mu_t(s) \sum_a \pi_t(a|s) \log \pi_t(a|s) && \text{Max. Caus. Ent} \\ \text{s.t. } \forall t : \quad & \sum_t \sum_s \mu_t(s) \phi(s) = \psi(\pi^*), \quad \sum_a \pi(a|s) = 1, \forall s && \text{Match Features} \\ & \mu_t(s') = \sum_{s,a} \mu_{t-1}(s) \pi_{t-1}(a|s) p(s'|s, a), && \text{State distribution consistency} \end{aligned}$$

State distribution at time step t has to be consistent with:

- state distribution and policy at time step $t-1$
- system dynamics

Maximum-Causal-entropy IRL [Ziebart 2010]



Solution: $\pi_t(a|s) \propto \exp(\mathbf{w}^T \phi(s) + \mathbb{E}[V_{t+1}(s')|s, a])$

$V_t(s)$ is again a Lagrangian multiplier

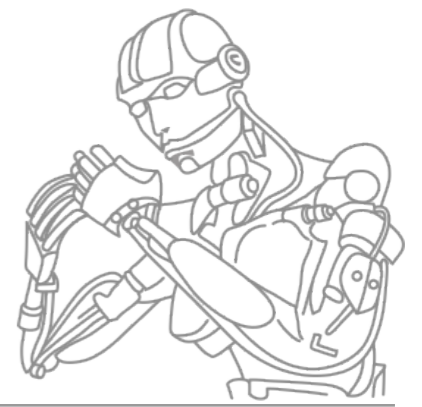
If we say $\mathbf{w}^T \phi(s)$ is the reward, this is a soft-max over the Q-function

$$\pi_t(a|s) \propto \exp(Q(s, a; \mathbf{w}))$$

This is still a convex problem:

Solution can be obtained by optimizing dual function

Can be done (relatively) efficiently



Outline of the Lecture

1. Introduction

- I. Comparison to Behavioral Cloning

2. Categories of IRL

- I. Maximum Margin
- II. Feature Matching by Max. Entropy
- III. Policy parametrized by rewards

3. Applications

4. Conclusion



Reward function parameterizing the policy class

Alternatively, we can assume the policy is soft-max in a Q-function

$$\pi_t(a|s; r, \alpha) = \frac{\exp(\alpha Q^*(s, a; r))}{\sum_a' \exp(\alpha Q^*(s, a'; r))}$$

where Q^* is the optimal Q-function for reward function $r(s)$, i.e.,

$$V^*(s; r) = \max_a Q^*(s, a; r)$$

$$Q^*(s, a; r) = r(s) + \mathbb{E}[V^*(s'; r)|s, a]$$

Then we can evaluate the likelihood of seeing a set of state-action pairs as follows:

$$\log p(D|r, \alpha) = \sum_i \log \pi(a_i|s_i; r, \alpha)$$

Is equivalent to “a smarter” Behavior Cloning!

Reward function parameterizing the policy class



Ziebart's approach can also be shown to be equivalent!

Can be extended to Bayesian setup:

Put prior on parameters of reward function

$$p(r, \alpha | D) = \frac{p(D | r, \alpha) p(r, \alpha)}{p(D)}$$

- Ramachandran and Amir, AAI2007: MCMC method to sample from this distribution
- Neu and Szepesvari, UAI2007: gradient method to optimize the likelihood [MAP]



Open Directions

- Open directions:
 - Active inverse RL,
 - Inverse RL with minmax control
 - Inverse RL with partial observability
 - Inverse RL with learning stages (rather than observing optimal policy)
 - Many more ...
- Are you interested? We may have an excellent thesis for you!

Outline of the Lecture



1. Introduction

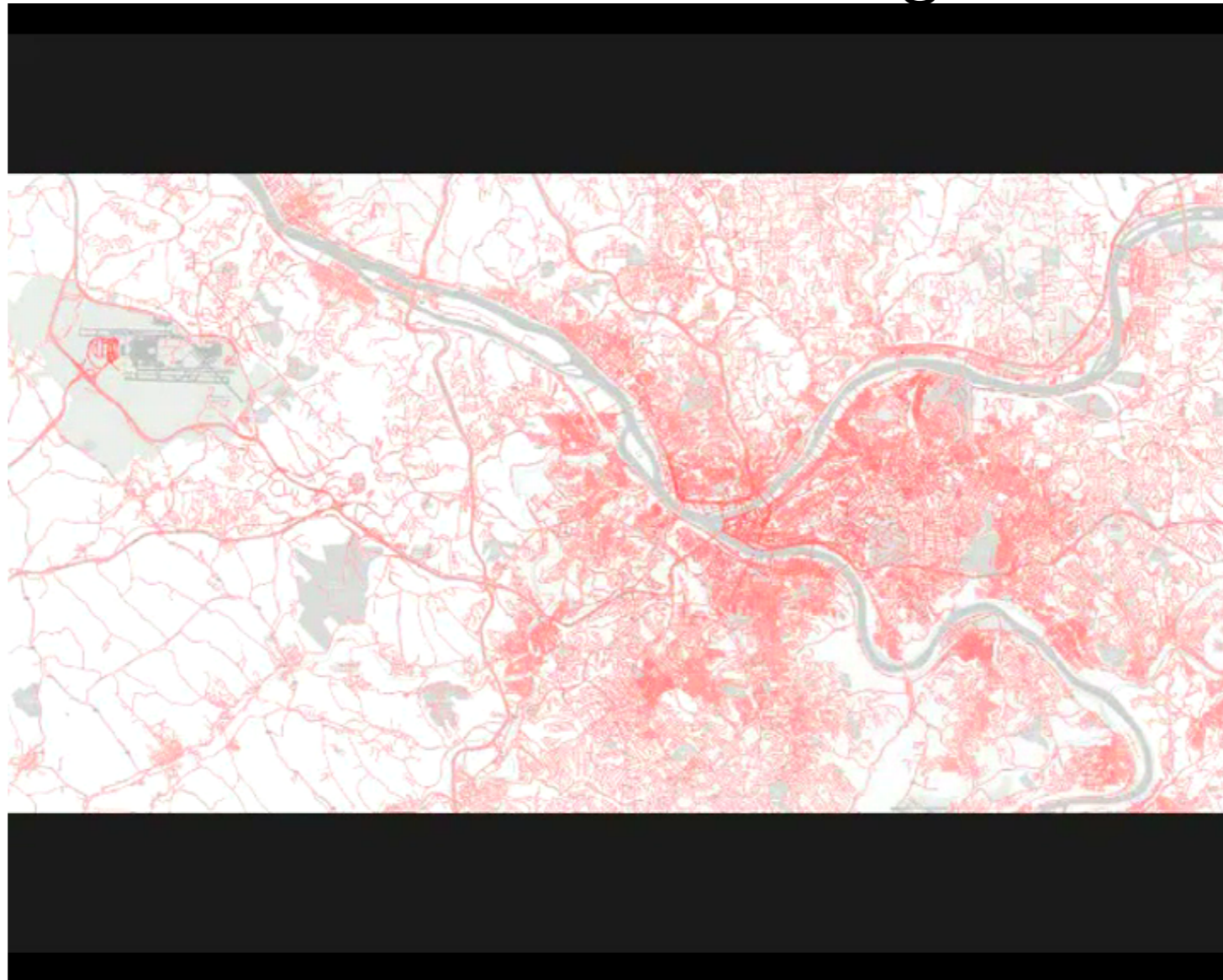
2. Contrast to behavioral cloning, historical sketch
Mathematical formulations for inverse RL
Example applications

▶ 3. Case studies

Urban navigation



- Reward function for urban navigation?



54 → destination prediction

- Ziebart, Maas, Bagnell and Dey AAI 2008



Outline of the Lecture

1. Introduction

- I. Comparison to Behavioral Cloning

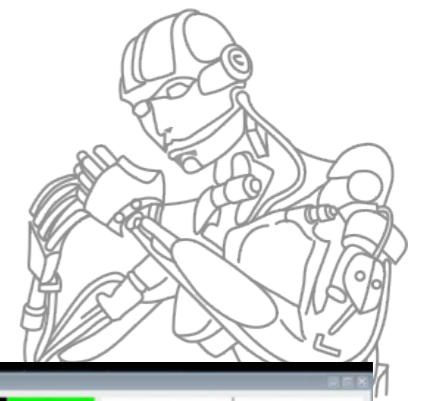
2. Categories of IRL

- I. Maximum Margin
- II. Feature Matching by Max. Entropy
- III. Policy parametrized by rewards

3. Applications

4. Conclusion

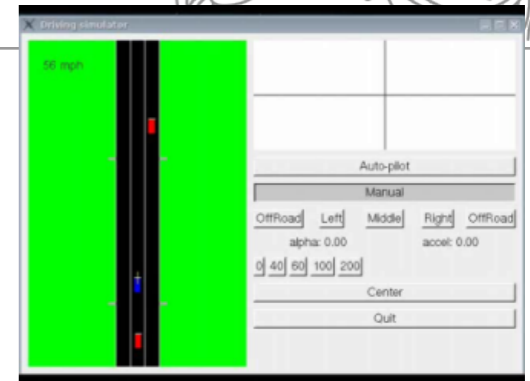
Examples



- Simulated highway driving

- Abbeel and Ng, ICML 2004,

- Syed and Schapire, NIPS 2007



- Aerial imagery based navigation

- Ratliff, Bagnell and Zinkevich, ICML 2006

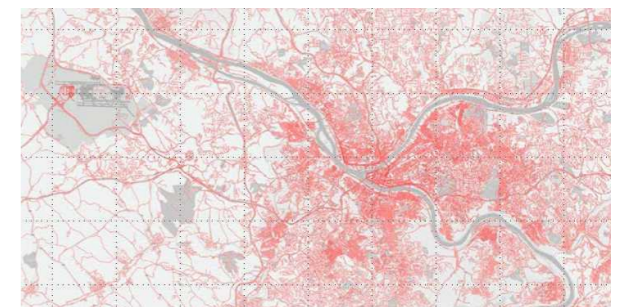


- Parking lot navigation

- Abbeel, Dolgov, Ng and Thrun, IROS 2008



- Urban navigation



56 Ziebart, Maas, Bagnell and Dey, AAAI 2008

Examples (ctd)



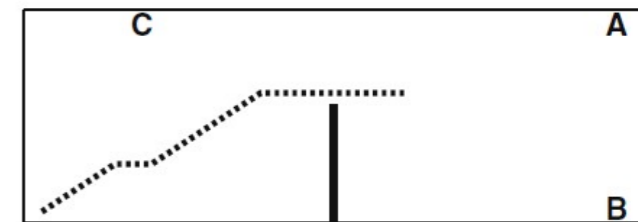
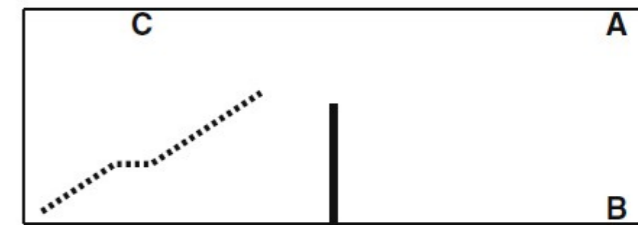
- Human path planning

- Mombaur, Truong and Laumond, AURO 2009



- Human goal inference

- Baker, Saxe and Tenenbaum, Cognition 2009



- Quadruped locomotion

- Ratliff, Bradley, Bagnell and Chestnutt, NIPS 2007

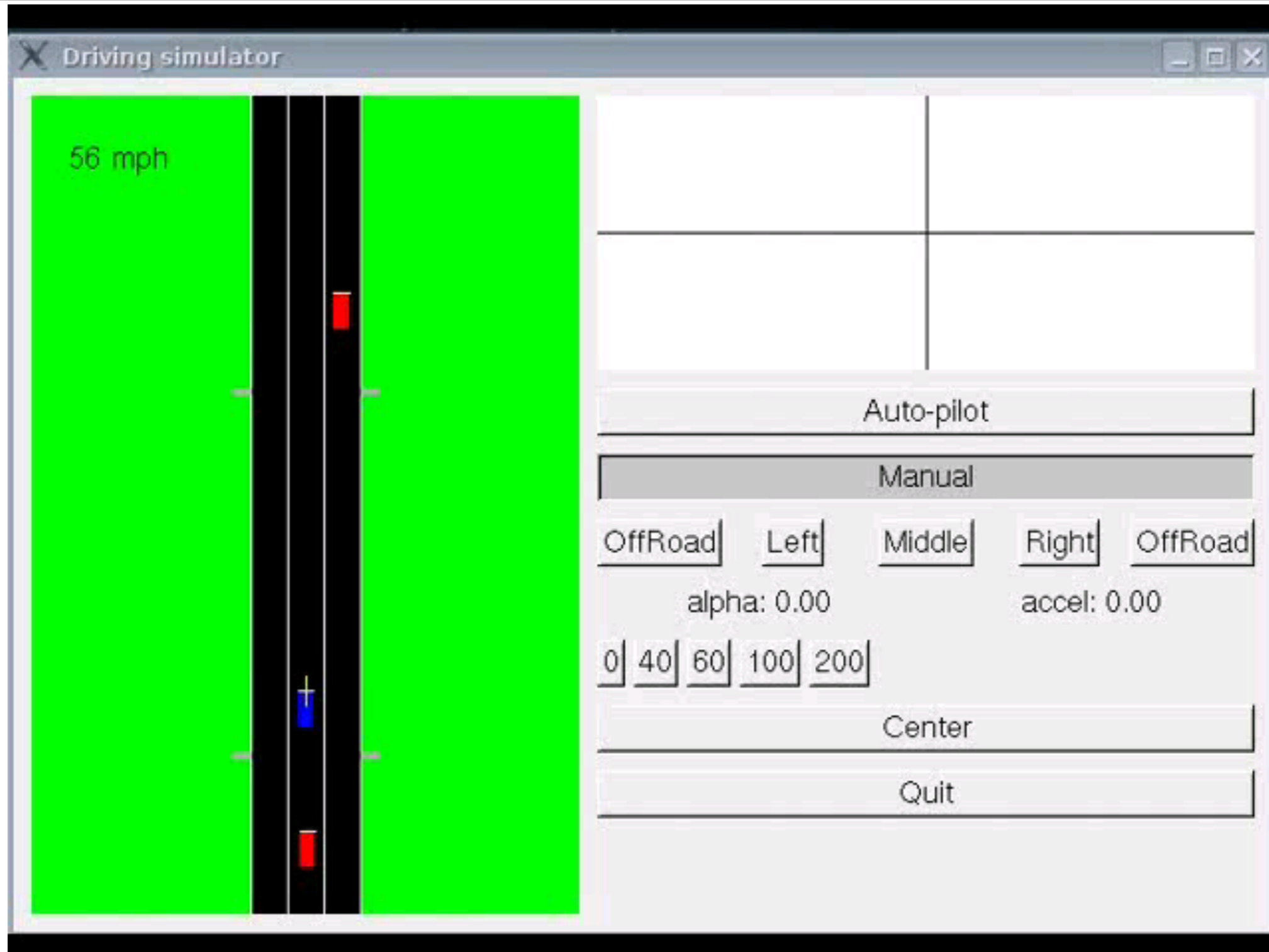




Lecture outline

- Case studies:
 - (1) Highway driving,
 - (2) Crusher,
 - (3) Parking lot navigation,
 - (4) Route inference,
 - (5) Human path planning,
 - (6) Human inverse planning,
 - (7) Quadruped locomotion
 - (8) Helicopter Acrobatics

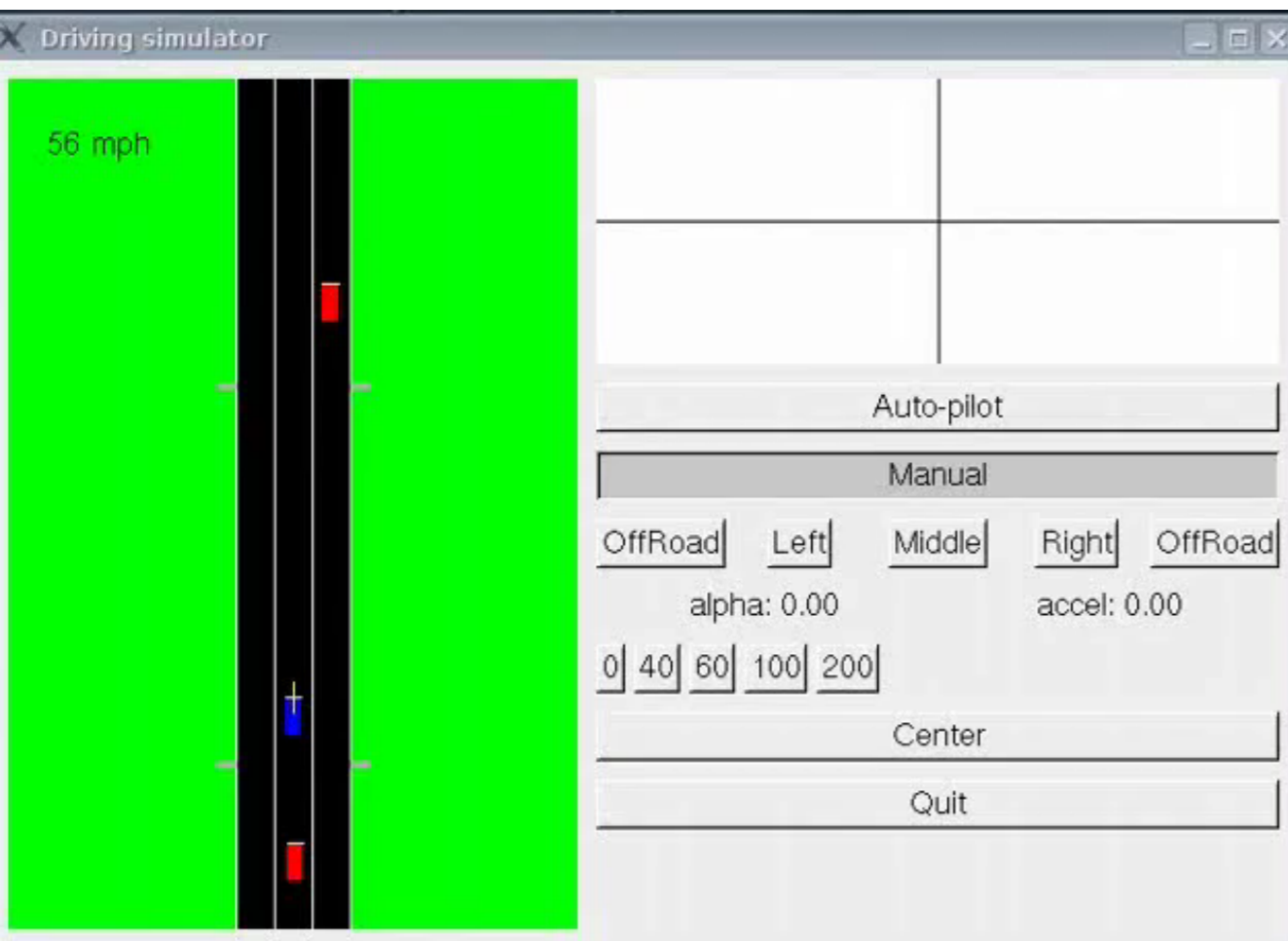
Simulated highway driving



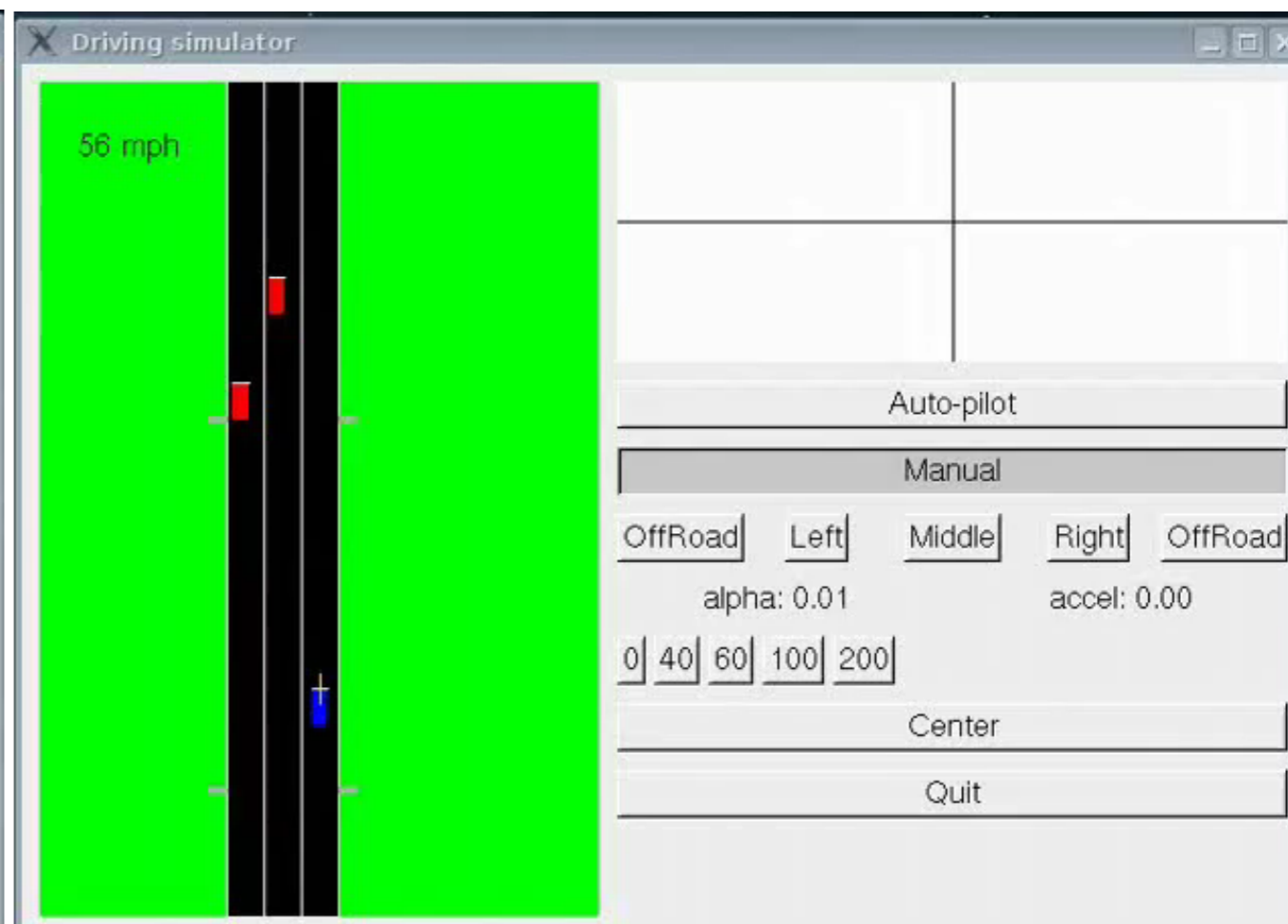
Highway driving



Teacher in Training World



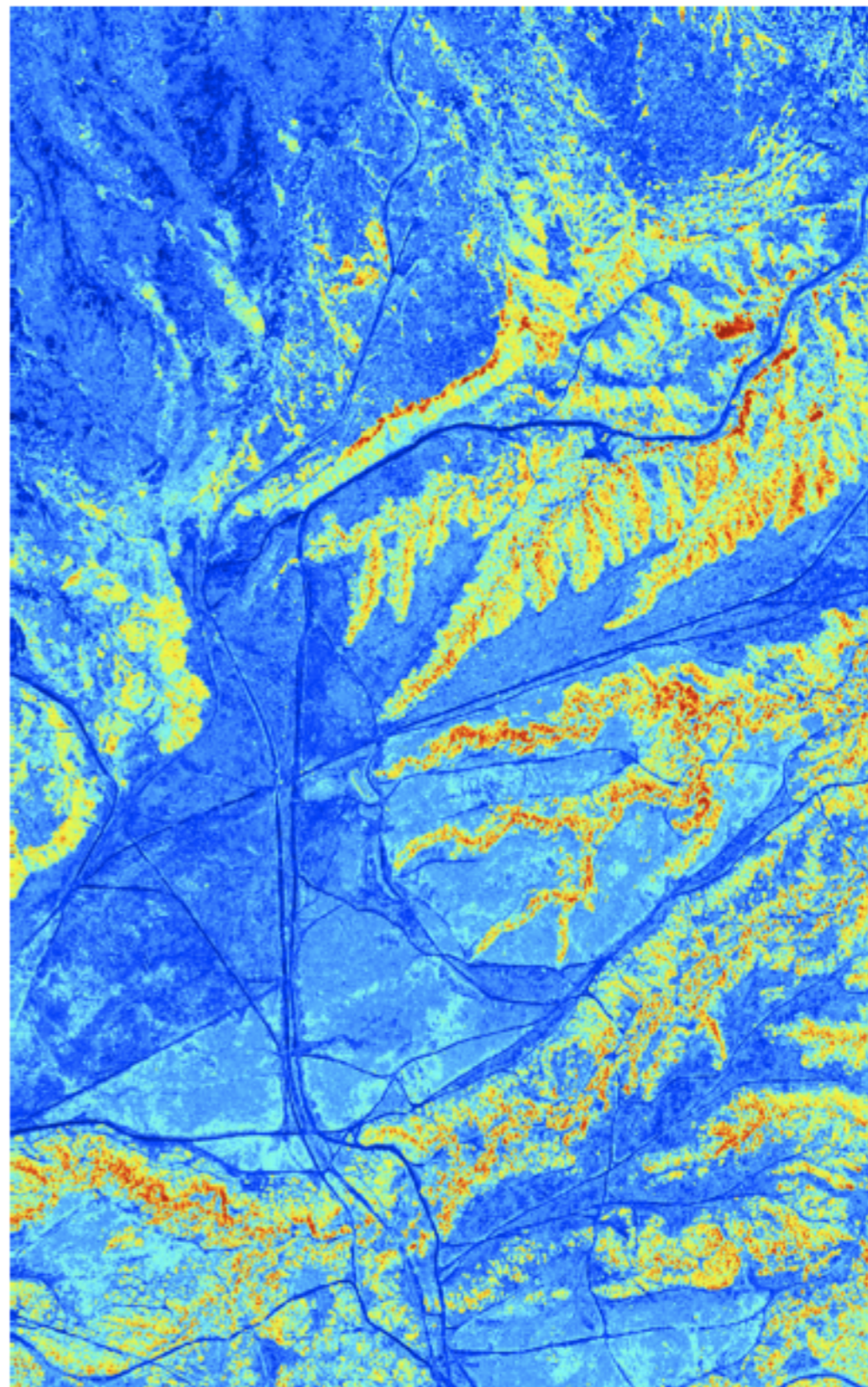
Learned Policy in Testing World



• Input:

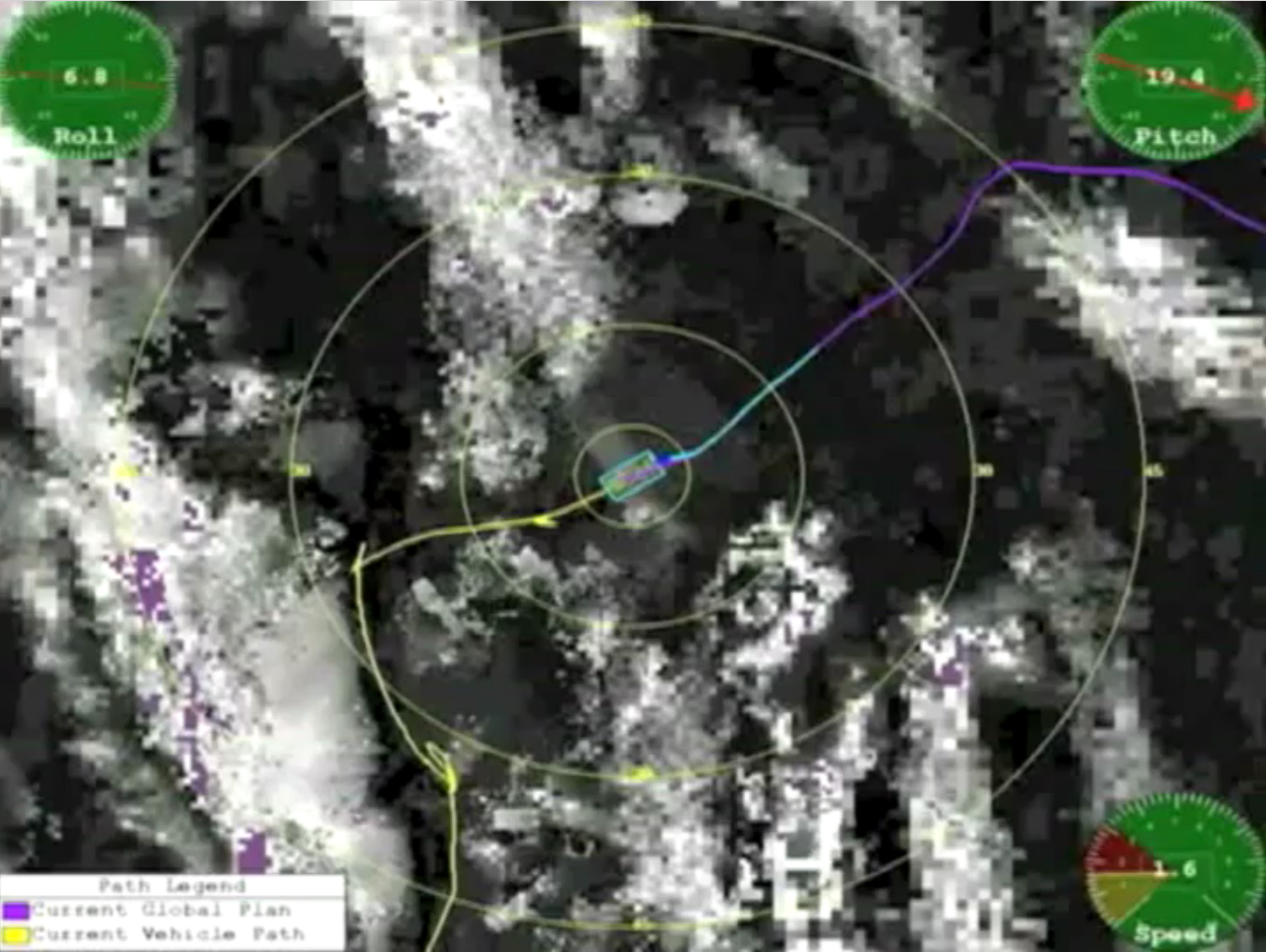
- Dynamics model / Simulator $P_{xu}(x_{t+1} | x_t, u_t)$ [Abbeel and Ng 2004]
- Teacher's demonstration: 1 minute in "training world"
- Note: R^* is unknown.
- Reward features: 5 features corresponding to lanes/shoulders; 10 features corresponding to presence of other car in current lane at different distances

Max margin



[Ratliff + al, 2006/7/8]

Max-margin





Parking lot navigation



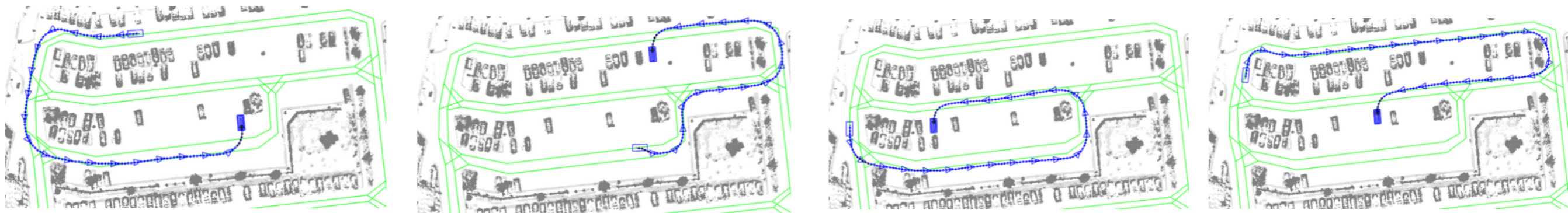
- Reward function trades off:
 - Staying “on-road,”
 - Forward vs. reverse driving,
 - Amount of switching between forward and reverse,
 - Lane keeping,
 - On-road vs. off-road,
 - Curvature of paths.

[Abbeel et al., IROS 08]

Experimental setup



- Demonstrate parking lot navigation on “train parking lots.”



- Run our apprenticeship learning algorithm to find the reward function.
- Receive “test parking lot” map + starting point and destination.
- Find the trajectory that maximizes the **learned reward function** for navigating the test parking lot.

Nice driving style

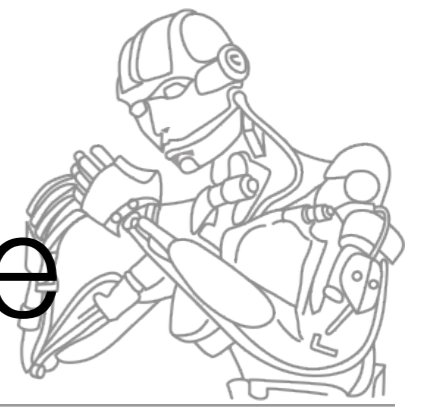


Sloppy driving-style



QuickTime™ and a
JVT/AVC Coding decompressor
are needed to see this picture.

“Don’t mind reverse” driving-style



QuickTime™ and a
JVT/AVC Coding decompressor
are needed to see this picture.



**Only 35% of routes are
“fastest”** (Letchner, Krumm, & Horvitz
2006)

Time



Money

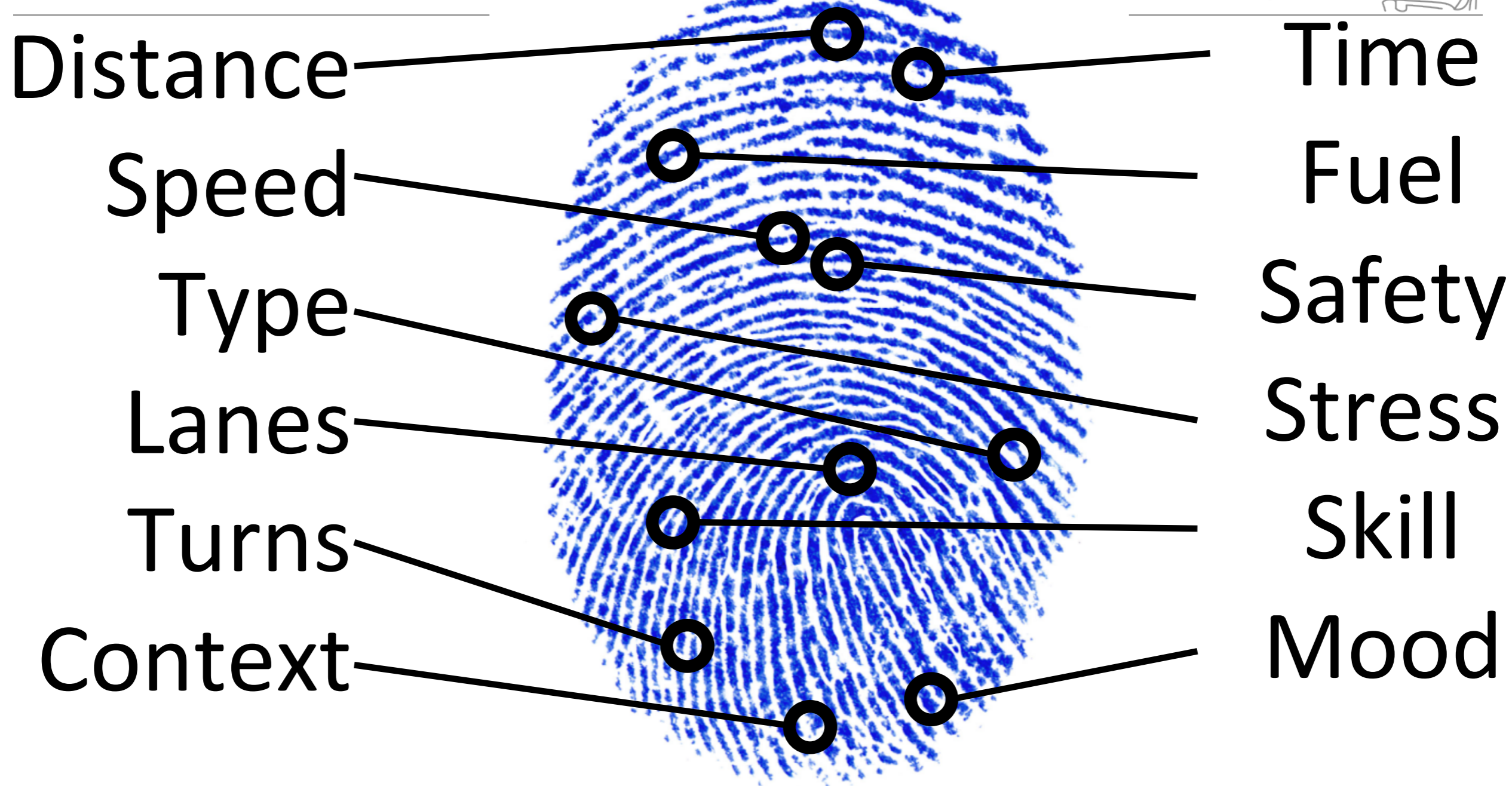


Stress



Skill





Data Collection

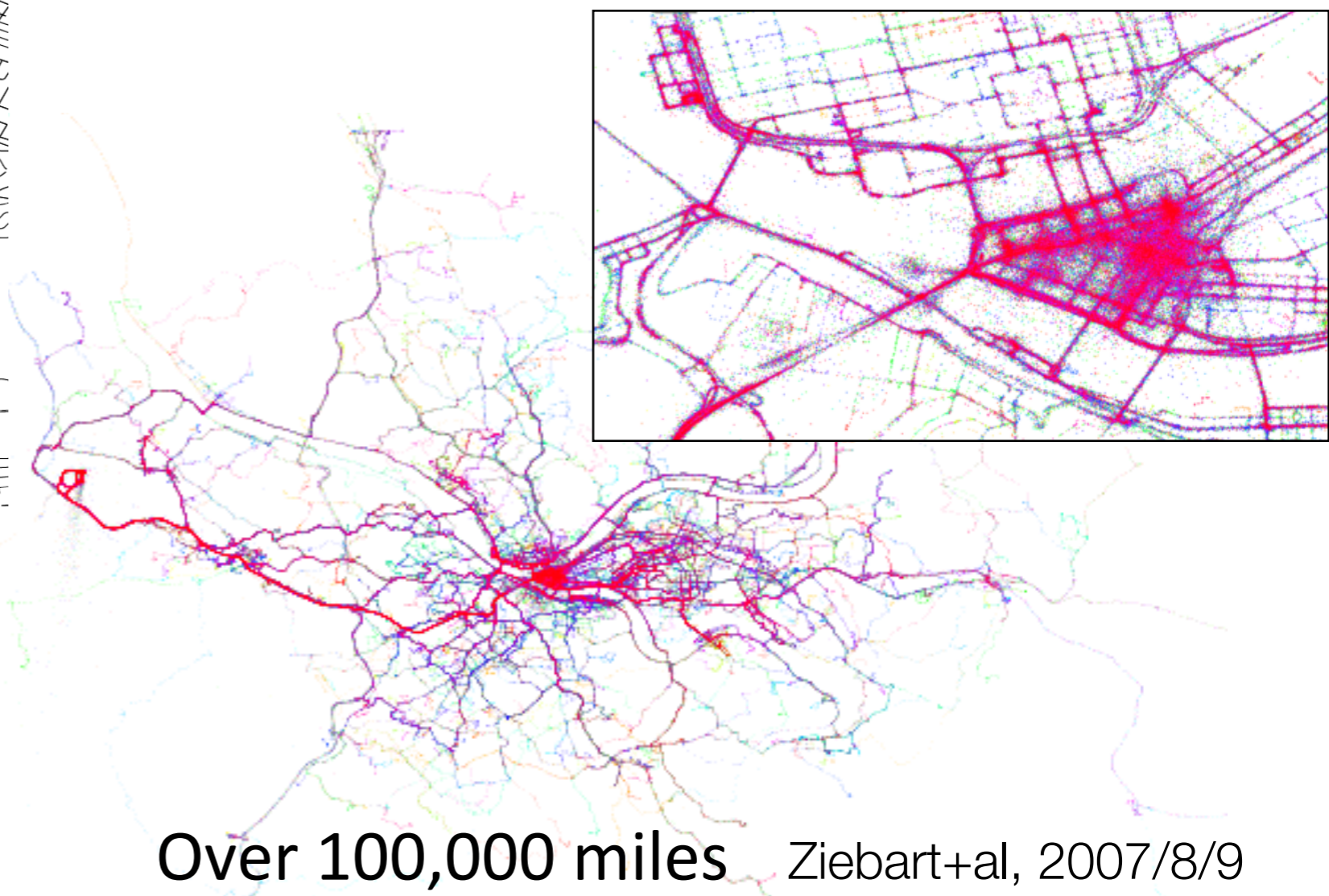


25 Taxi Drivers



Length
Speed
Road Type
Lanes

Accidents
Construction
Congestion
Time of day



Over 100,000 miles Ziebart+al, 2007/8/9

Destination Prediction



Equivalent: Pedestrian Prediction



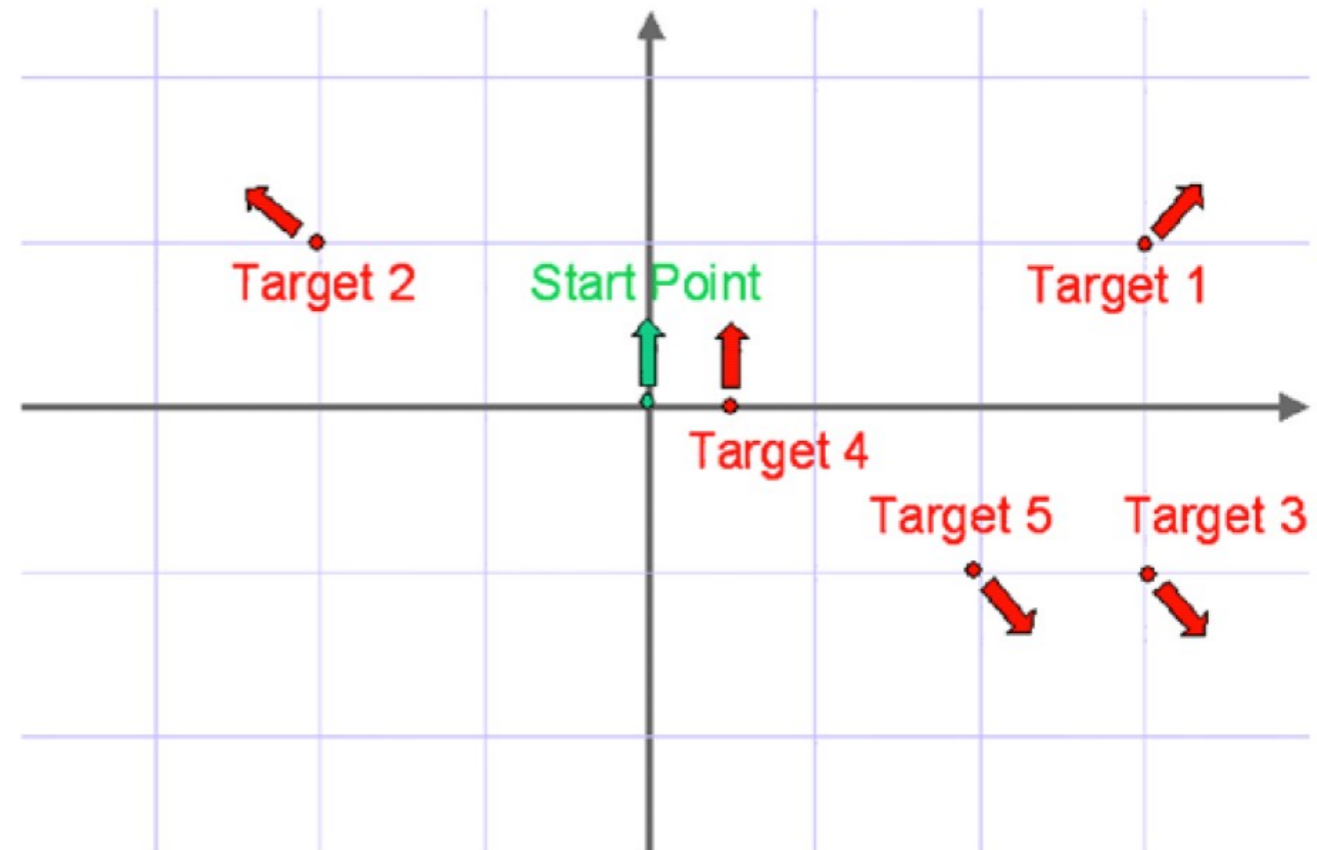
QuickTime™ and a
JVT/AVC Coding decompressor
are needed to see this picture.

Human path planning

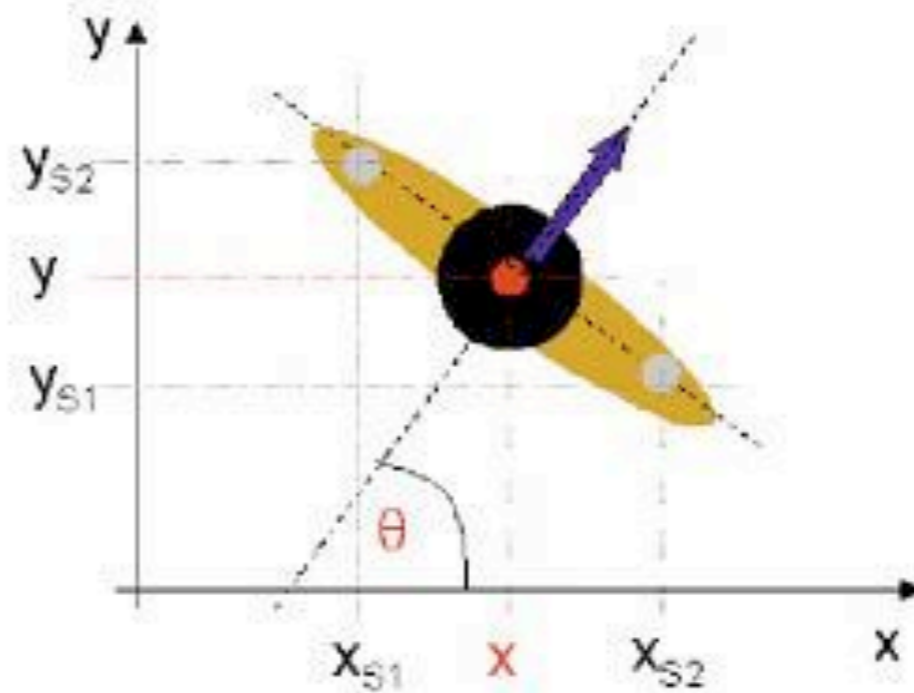
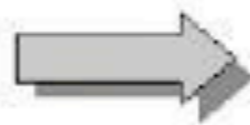
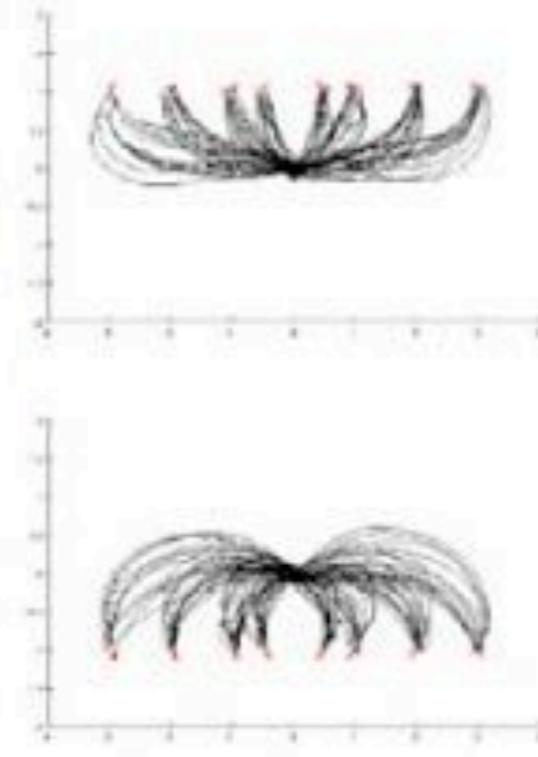


■ Reward features:

- Time to destination
- $(\text{Forward acceleration})^2$
- $(\text{Sideways acceleration})^2$
- $(\text{Rotational acceleration})^2$
- Integral $(\text{angular error})^2$



Experimental Setup



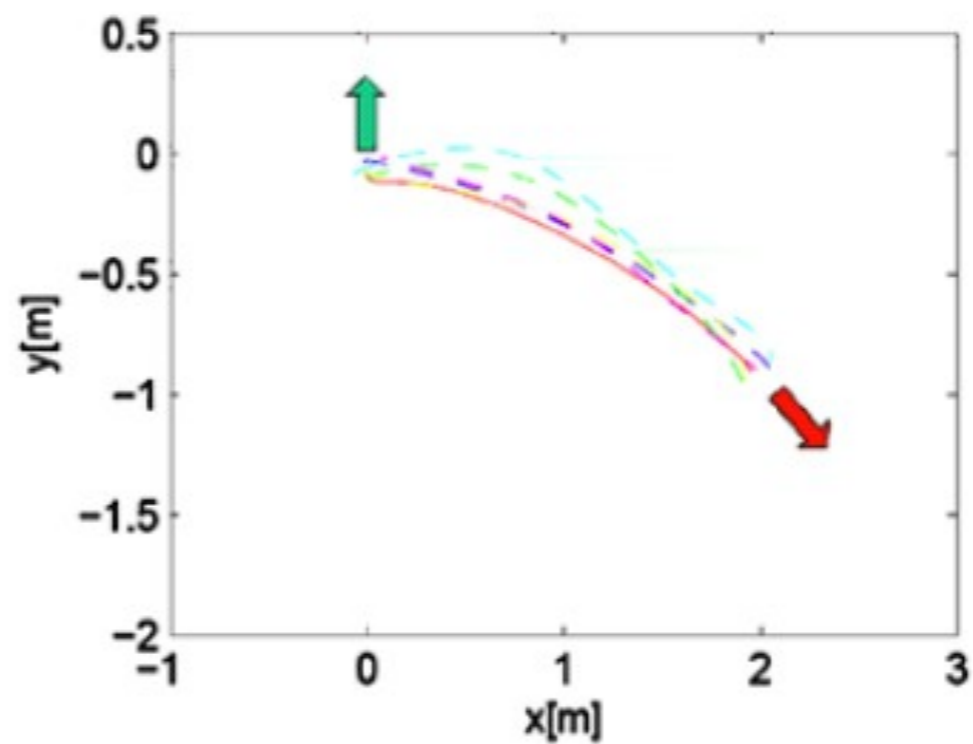
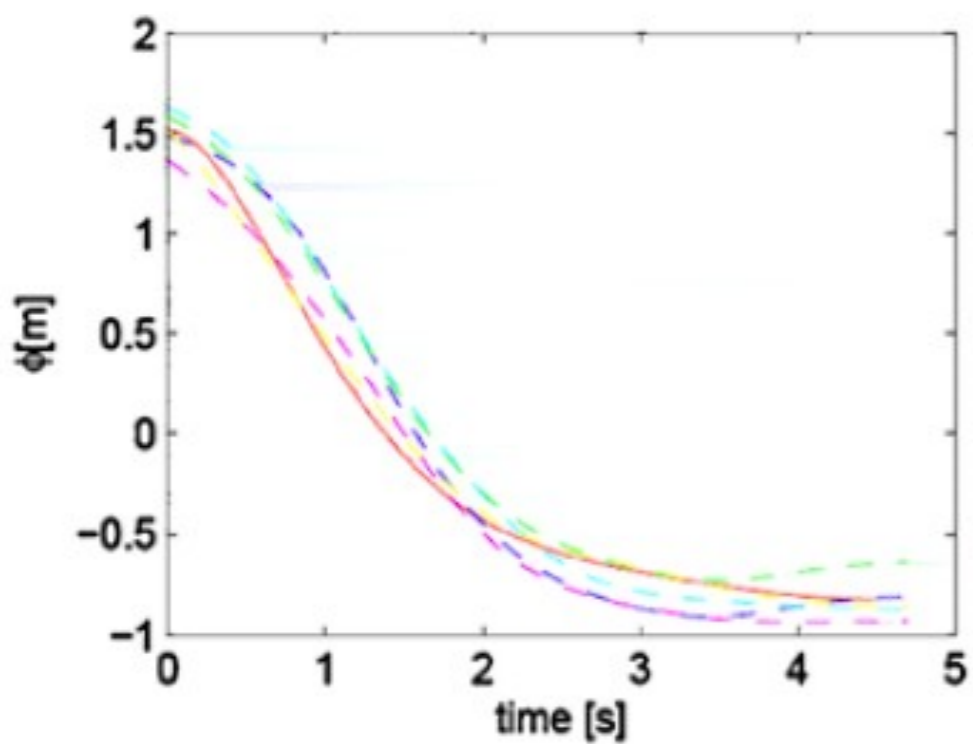
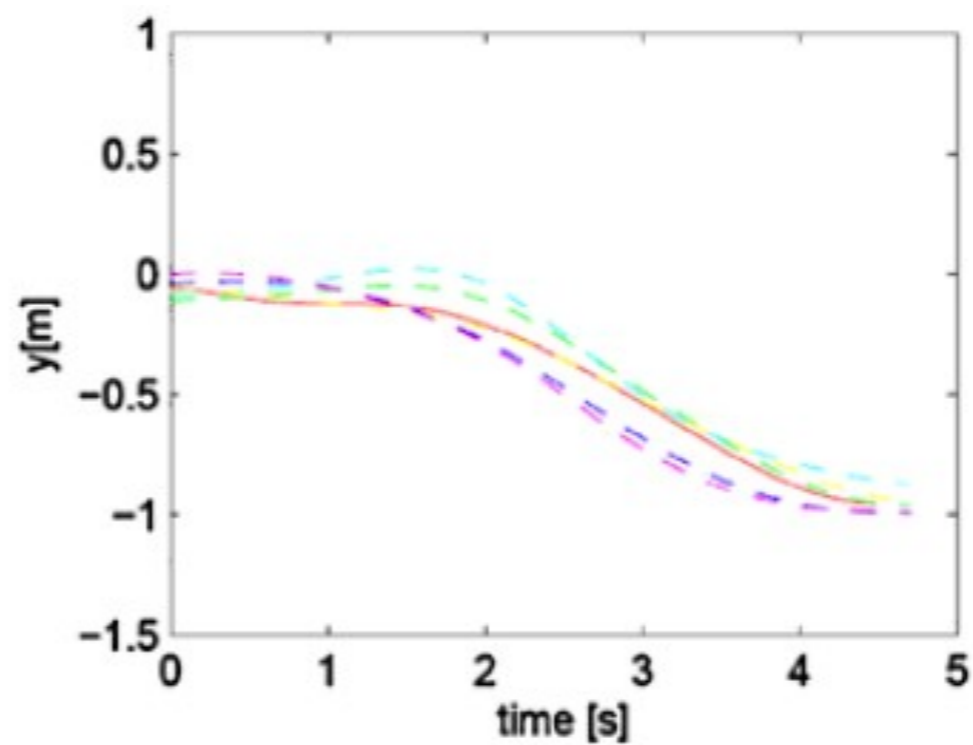
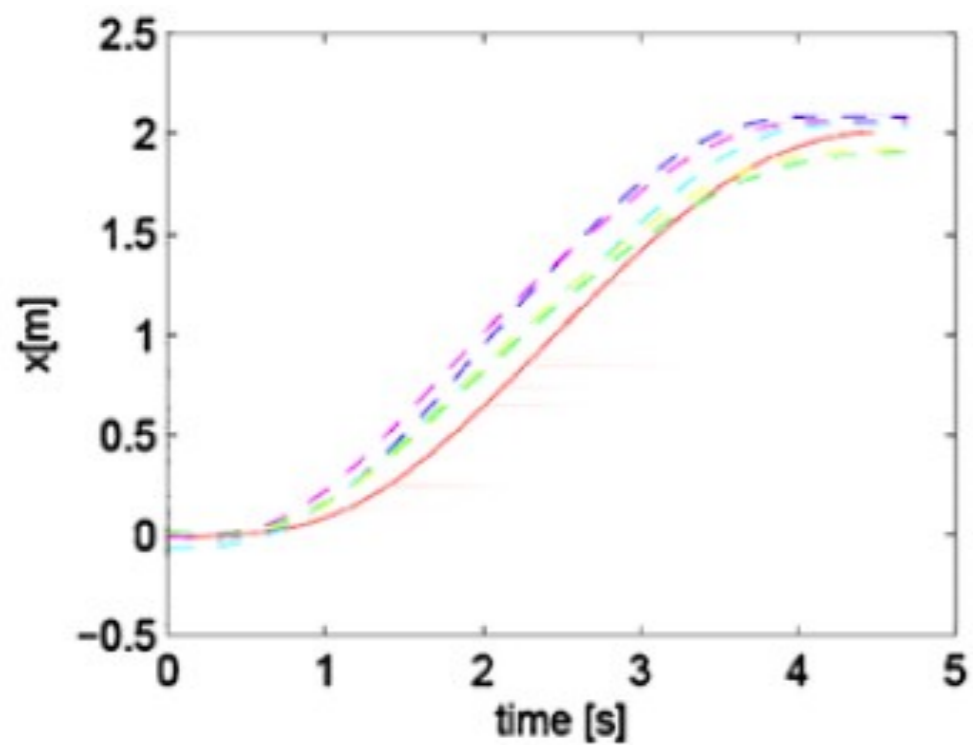
Human path planning



■ Result:

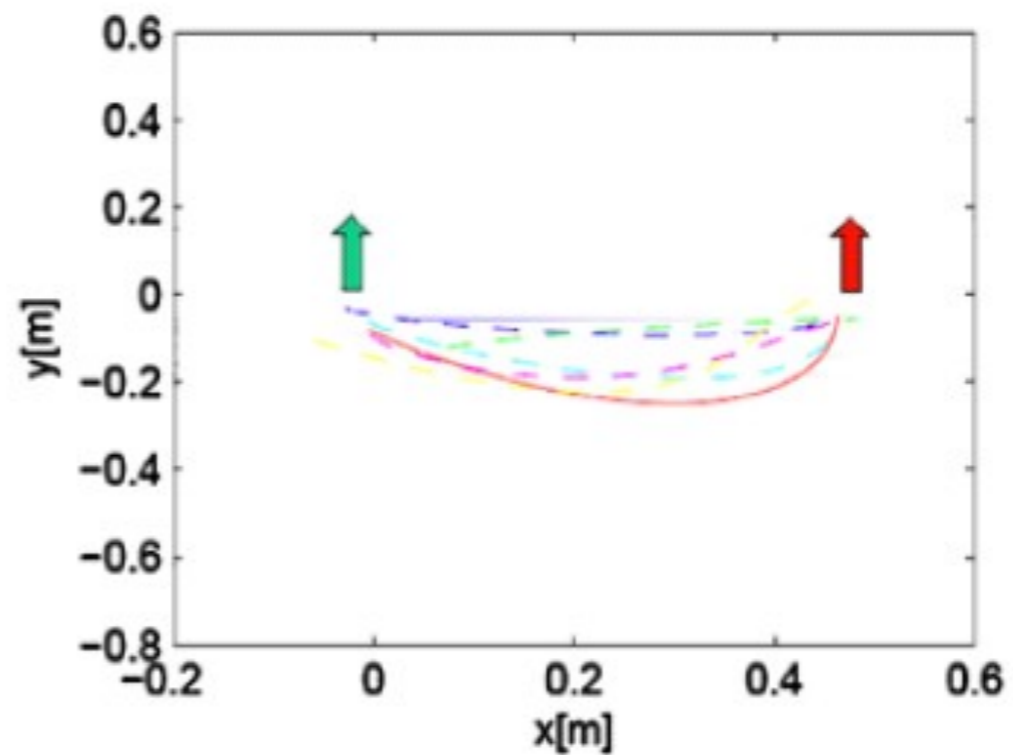
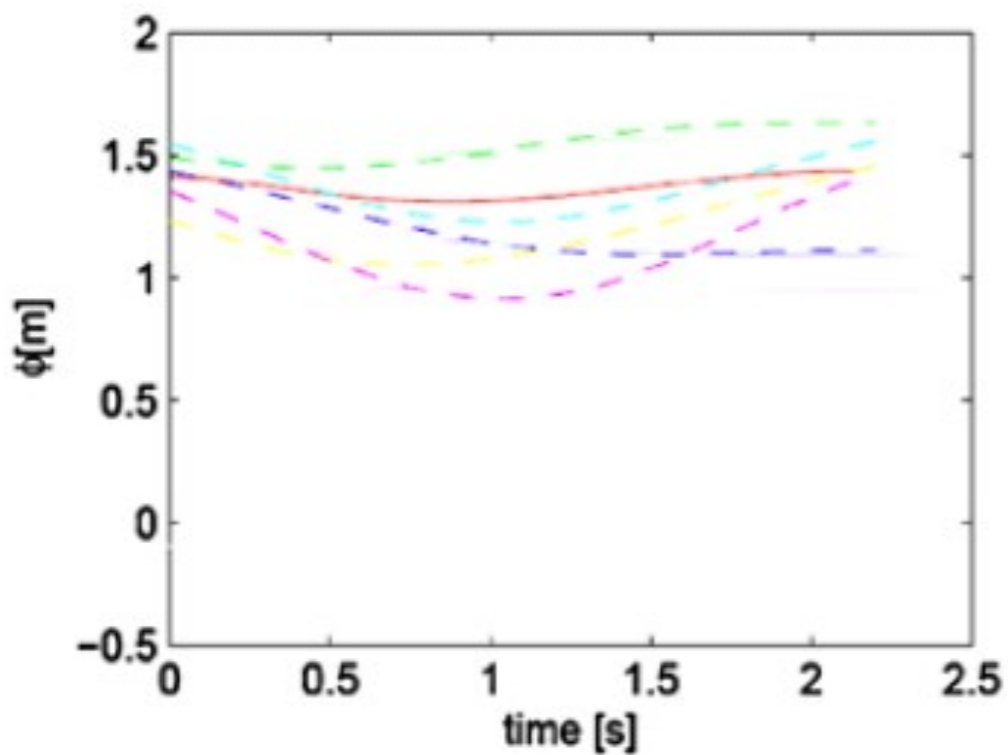
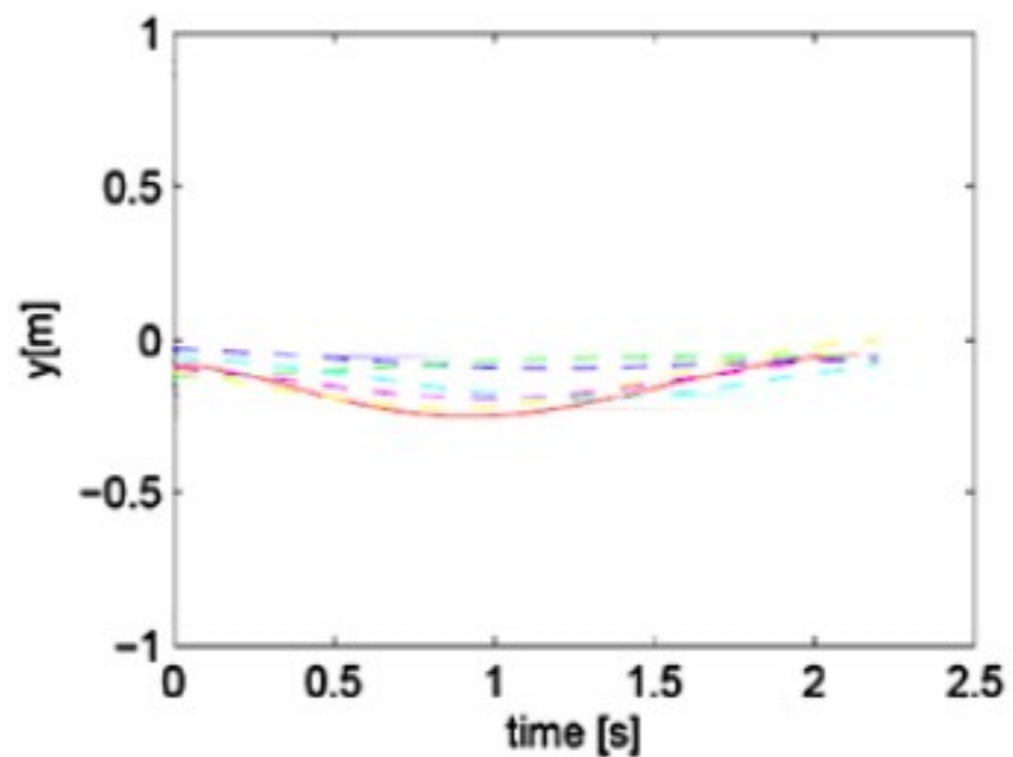
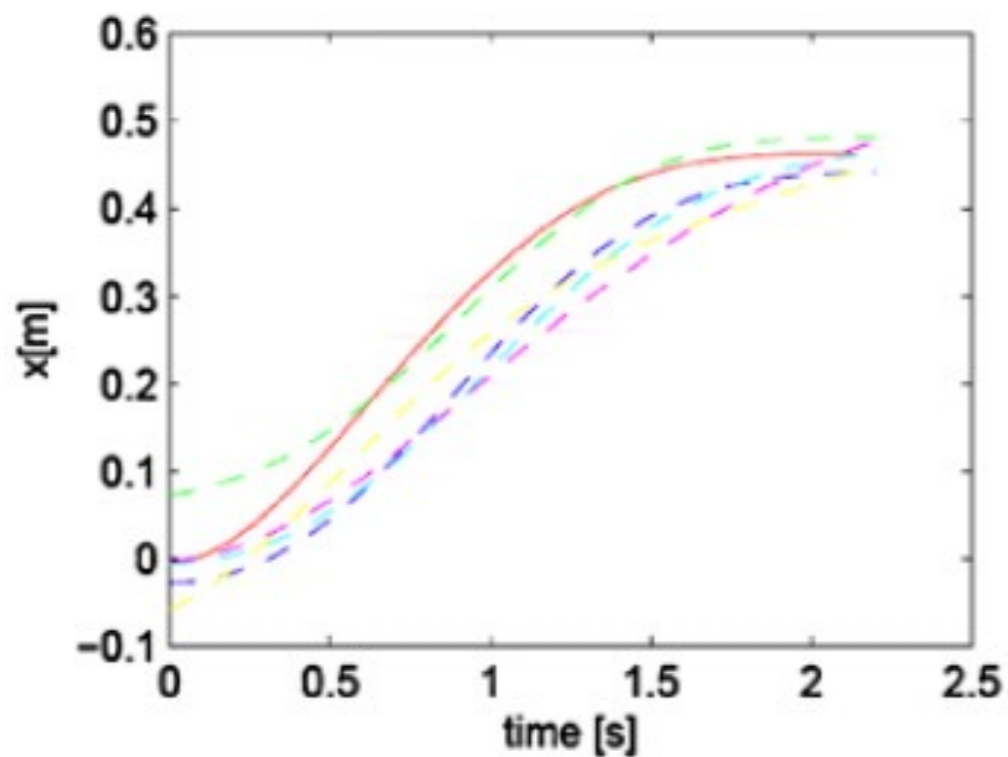
- Time to destination: 1
- (Forward acceleration)² 1.2
- (Sideways acceleration)² 1.7
- (Rotational acceleration)² 0.7
- Integral (angular error)² 5.2

Human path planning



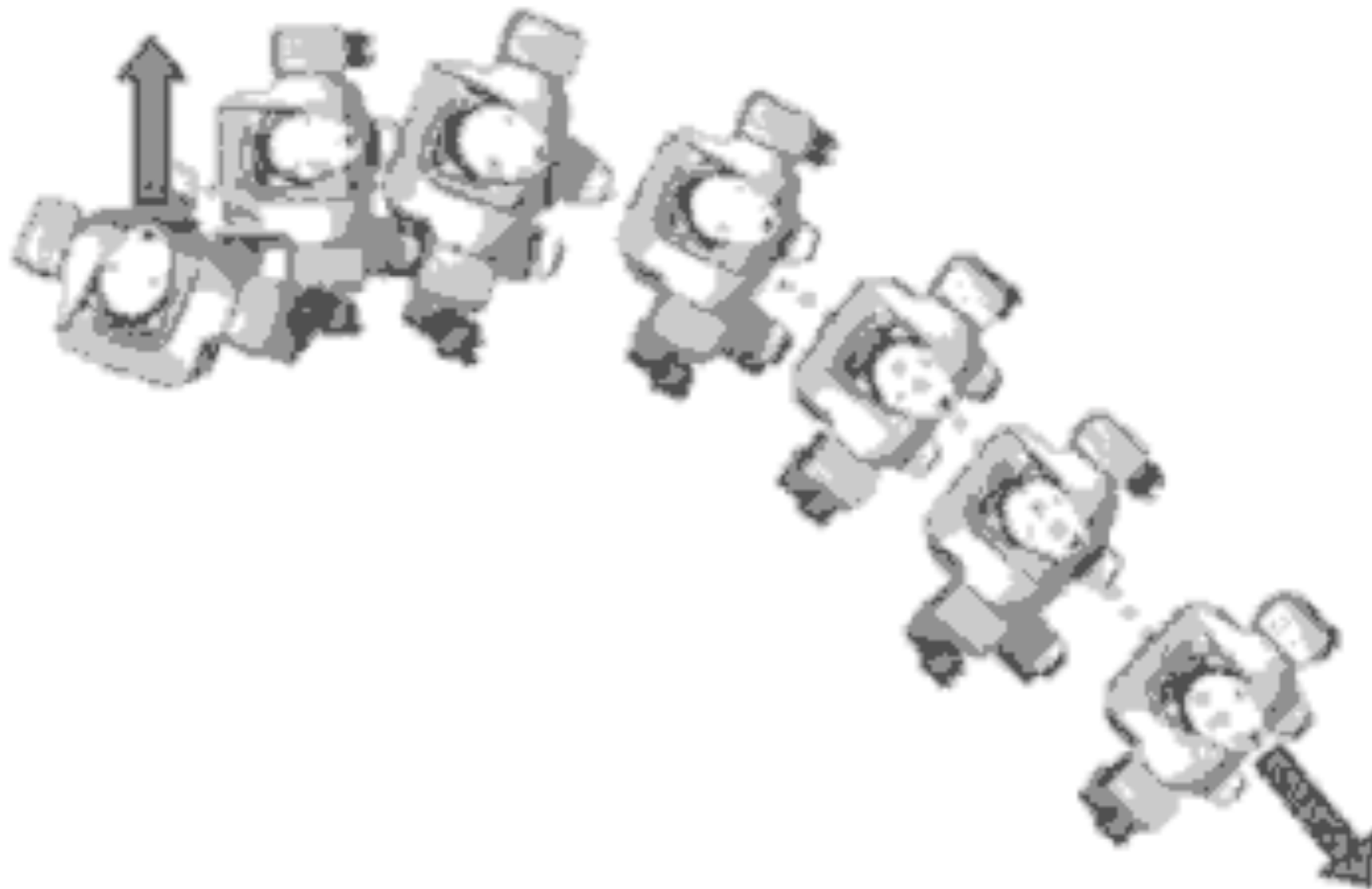
[Mombaur, Truong, Laumond, 2009]

Human path planning



[Mombaur, Truong, Laumond, 2009]

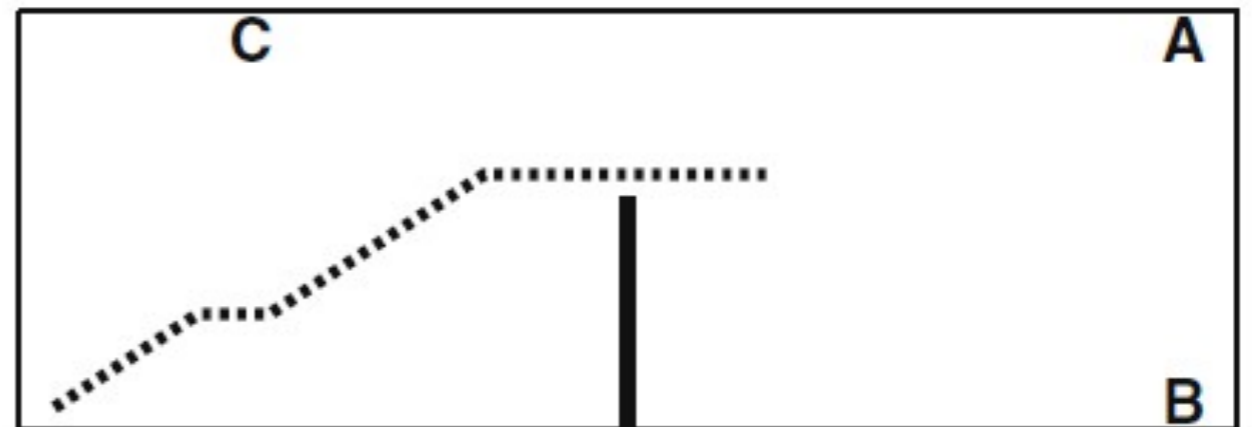
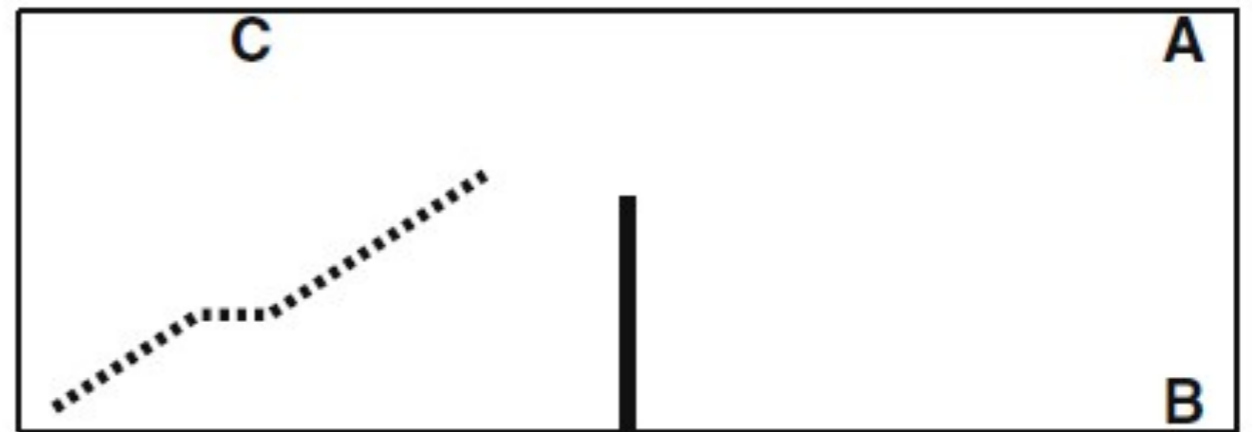
Transfer to a Humanoid



Goal inference



- Observe partial paths, predict goal. Goal could be either A, B, or C.
- + HMM-like extension: goal can change (with some probability over time).

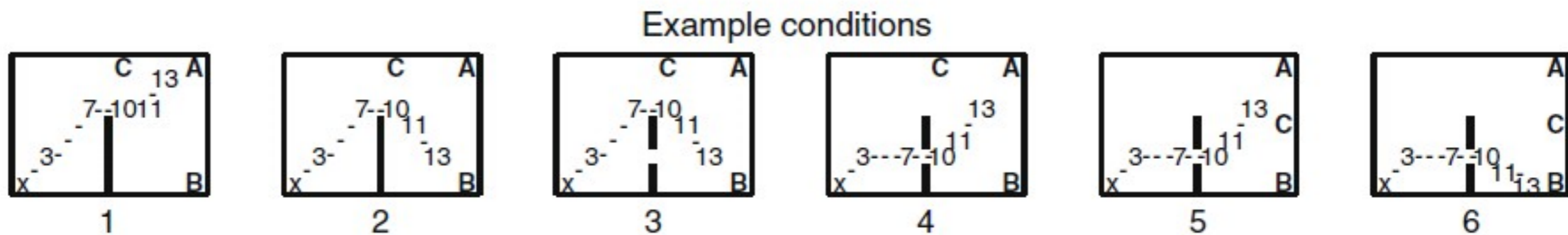


[Baker, Saxe, Tenenbaum, 2009]

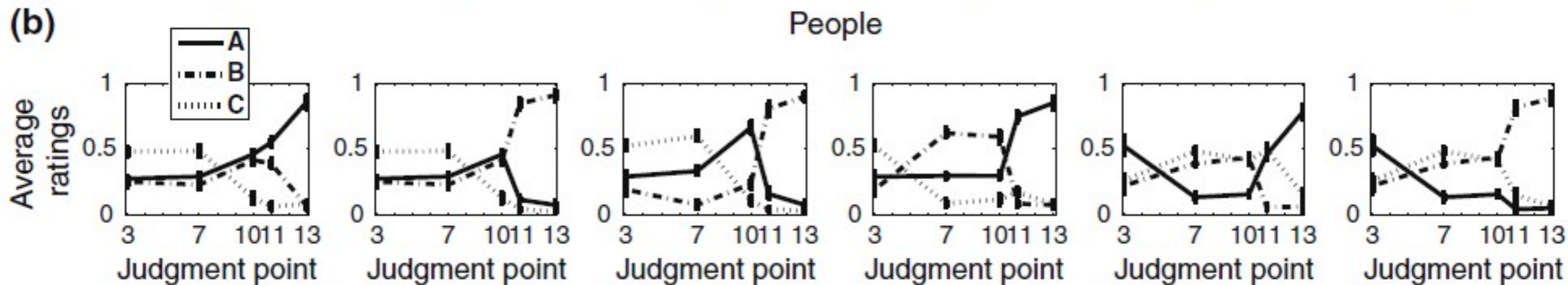
Goal inference



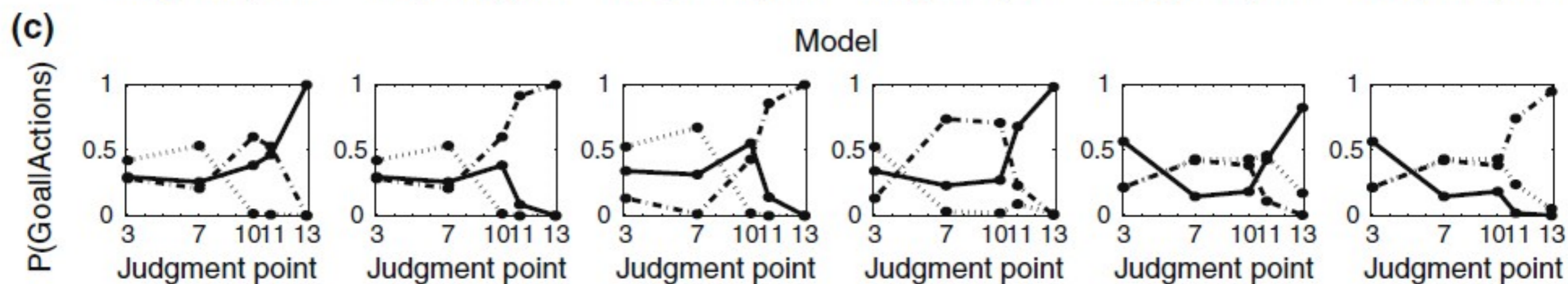
(a)



(b)



(c)



Quadruped

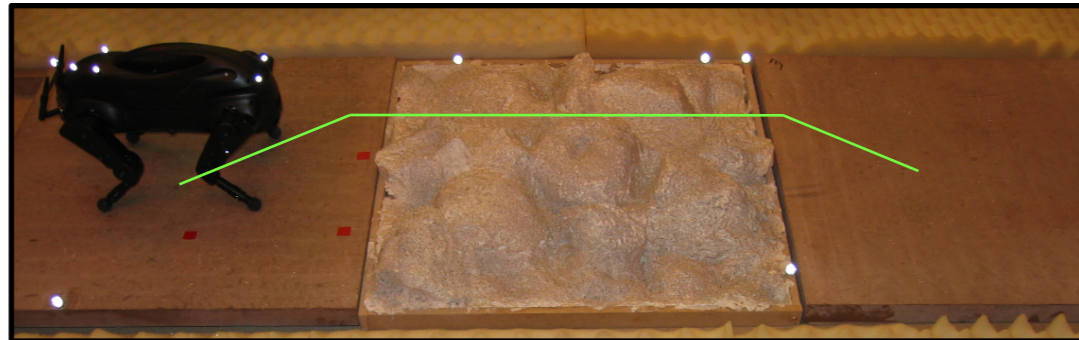


- Reward function trades off 25 features.

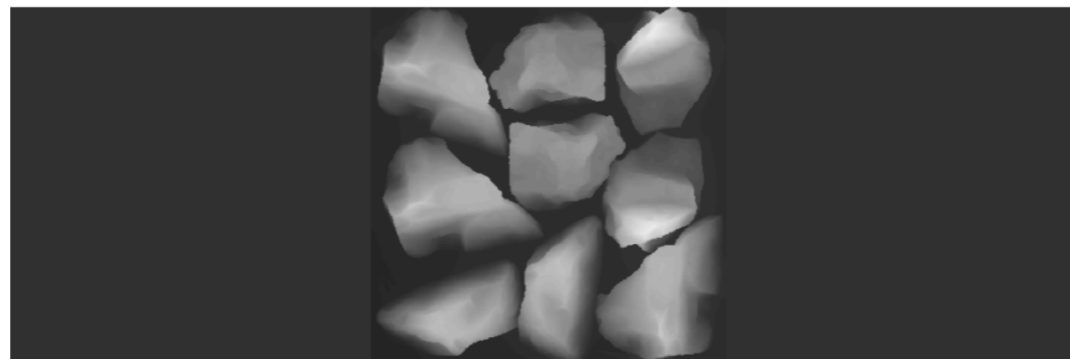
Experimental setup



- Demonstrate path across the “training terrain”



- Run the apprenticeship learning algorithm to find the reward function
- Receive “testing terrain” --- height map.



- Find the optimal policy with respect to the **learned reward function** for crossing the testing terrain.

Little Dog: CMU Team



Ratliff + al, 2007

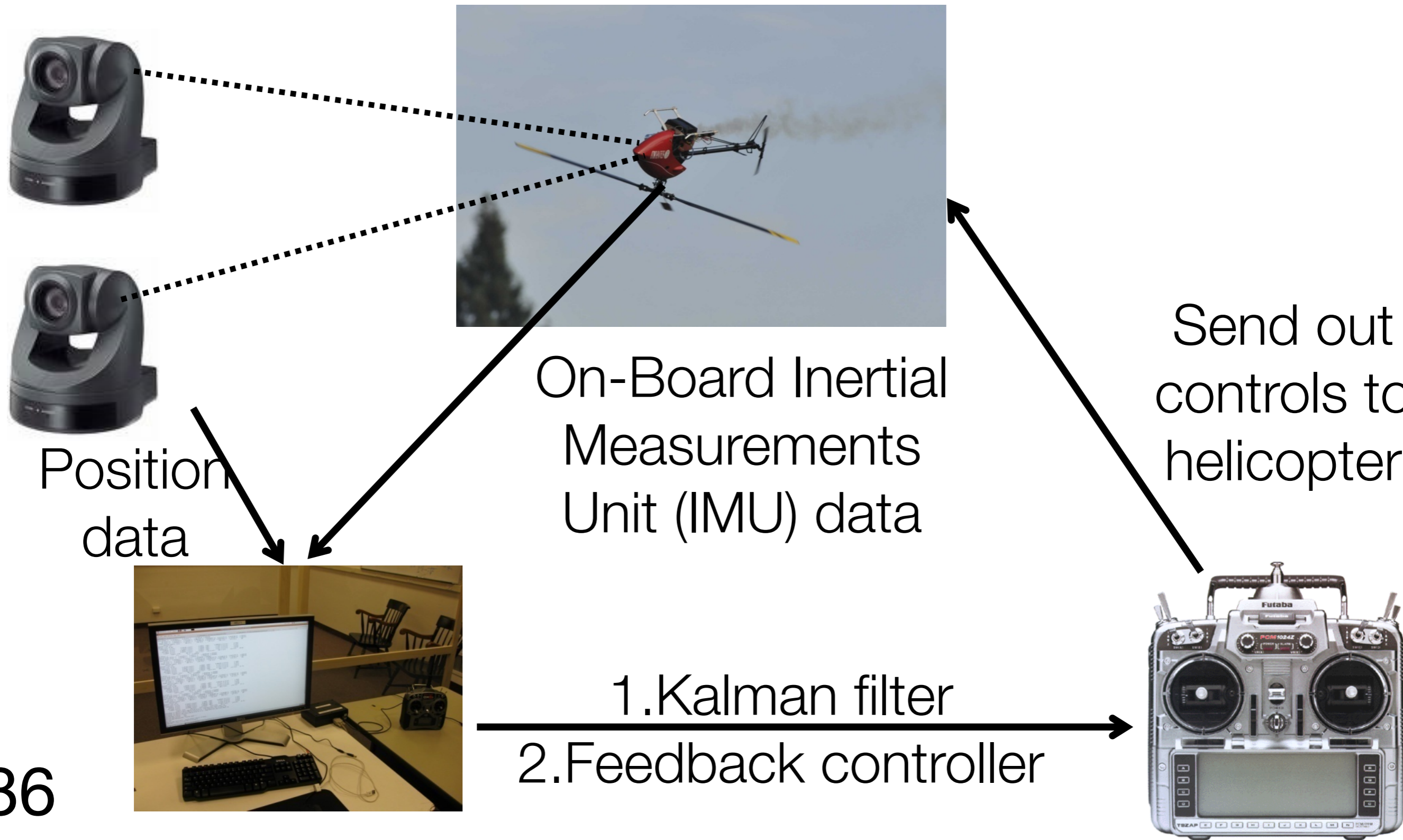
Remainder of lecture: extreme helicopter flight



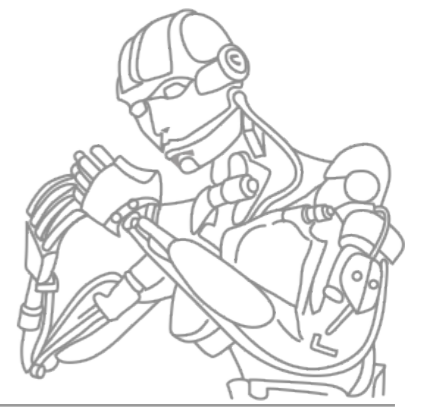
- How does helicopter dynamics work
- Autonomous helicopter setup
- Application of inverse RL to autonomous helicopter flight



Autonomous helicopter setup



Helicopter dynamics



- 4 control inputs:
 - Main rotor collective pitch
 - Main rotor cyclic pitch (roll and pitch)
 - Tail rotor collective pitch



Experimental setup for the helicopter



1. Our expert pilot demonstrates the airshow several times.



2. Learn (by solving a joint optimization problem):

- Reward function---trajectory.
- Dynamics model---trajectory-specific local model.

3. Fly autonomously:

- Inertial sensing + vision-based position sensing → (extended) Kalman filter
- Receding horizon differential dynamic programming (DDP) feedback controller (20Hz)
- **Learning to fly new aerobatics takes < 1 hour**

Results!

Summary



What you should know:

- ➔ Why is inverse RL useful / better than direct imitation learning?
- ➔ Algorithmic Challenges in IRL
- ➔ Different methods that use IRL, all are linear in features
- ➔ Why maximum margin?
- ➔ Why max. entropy?