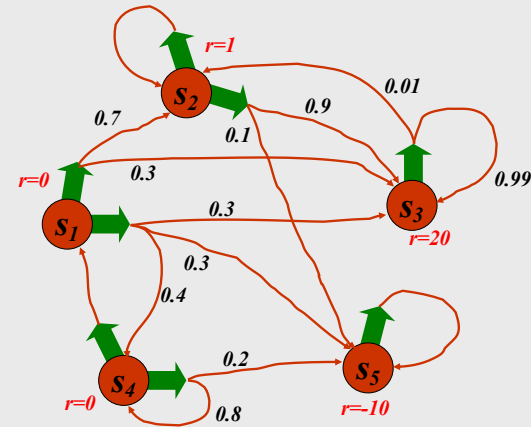# CSE-571
# Robotics

**Planning and Control:**

**Markov Decision Processes**

---

## Markov Decision Process (MDP)

---

## Markov Decision Process (MDP)

- **Given:**
- States $x$
- Actions $u$
- Transition probabilities $p(x'|u,x)$
- Reward / payoff function $r(x,u)$

- **Wanted:**
- Policy $\pi(x)$ that maximizes the future expected reward

---

## Rewards and Policies

- Policy (general case):

$$\pi:\quad z_{1:t-1}, u_{1:t-1} \rightarrow u_t$$

- Policy (fully observable case):

$$\pi:\quad x_t \rightarrow u_t$$

- Expected cumulative payoff:

$$R_T = E\left[\sum_{\tau=1}^{T} \gamma^\tau r_{t+\tau}\right]$$

- T=1: greedy policy
- T>1: finite horizon case, typically no discount
- T=infty: infinite-horizon case, finite reward if discount < 1

## Policies contd.

- Expected cumulative payoff of policy:

$$R_T^\pi(x_t) = E\left[\sum_{\tau=1}^{T} \gamma^\tau r_{t+\tau} \mid u_{t+\tau} = \pi\left(z_{1:t+\tau-1} u_{1:t+\tau-1}\right)\right]$$

- Optimal policy:

$$\pi^* = \operatorname*{argmax}_\pi \ R_T^\pi(x_t)$$

- 1-step optimal policy:

$$\pi_1(x) = \operatorname*{argmax}_u \ r(x,u)$$

- Value function of 1-step optimal policy:

$$V_1(x) = \gamma \max_u r(x,u)$$

5

## 2-step Policies

- Optimal policy:

$$\pi_2(x) = \operatorname*{argmax}_u \ \left[r(x,u) + \int V_1(x')p(x'\mid u,x)\,dx'\right]$$

- Value function:

$$V_2(x) = \gamma \max_u \ \left[r(x,u) + \int V_1(x')p(x'\mid u,x)\,dx'\right]$$

6

## T-step Policies

- Optimal policy:

$$\pi_T(x) = \operatorname*{argmax}_u \ \left[r(x,u) + \int V_{T-1}(x')p(x'\mid u,x)\,dx'\right]$$

- Value function:

$$V_T(x) = \gamma \max_u \ \left[r(x,u) + \int V_{T-1}(x')p(x'\mid u,x)\,dx'\right]$$

7

## Infinite Horizon

- Optimal policy:

$$V_\infty(x) = \gamma \max_u \ \left[r(x,u) + \int V_\infty(x')p(x'\mid u,x)\,dx'\right]$$

- Bellman equation

- Fix point is optimal policy

- Necessary and sufficient condition

8

## Value Iteration

- for all $x$ do
  $$\hat{V}(x) \leftarrow r_{min}$$
- endfor

- repeat until convergence
  - for all $x$ do
    $$\hat{V}(x) \leftarrow \gamma \max_{u} \left[ r(x,u) + \int \hat{V}(x')p(x'|u,x)dx' \right]$$
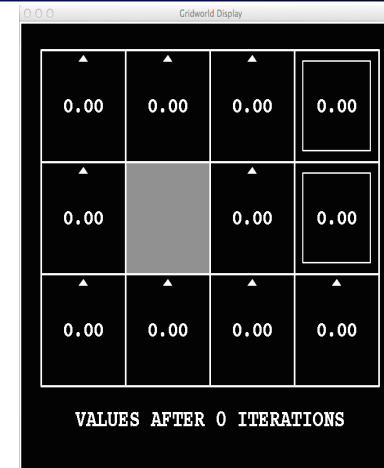  - endfor
- endrepeat
  $$\pi(x) = \arg\max_{u} \left[ r(x,u) + \int \hat{V}(x')p(x'|u,x)dx' \right]$$
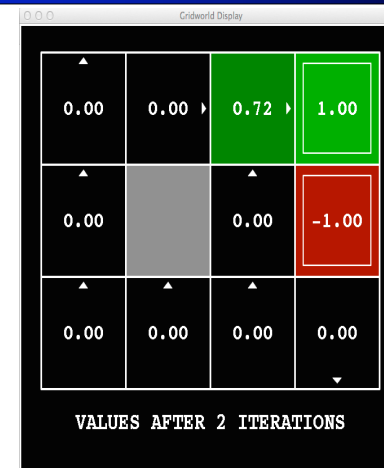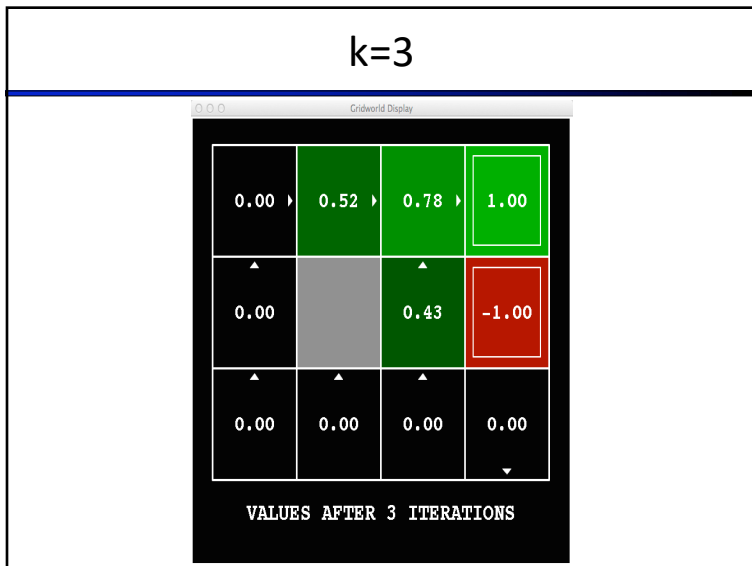
9

## k=0



Noise = 0.2
Discount = 0.9
Living reward = 0
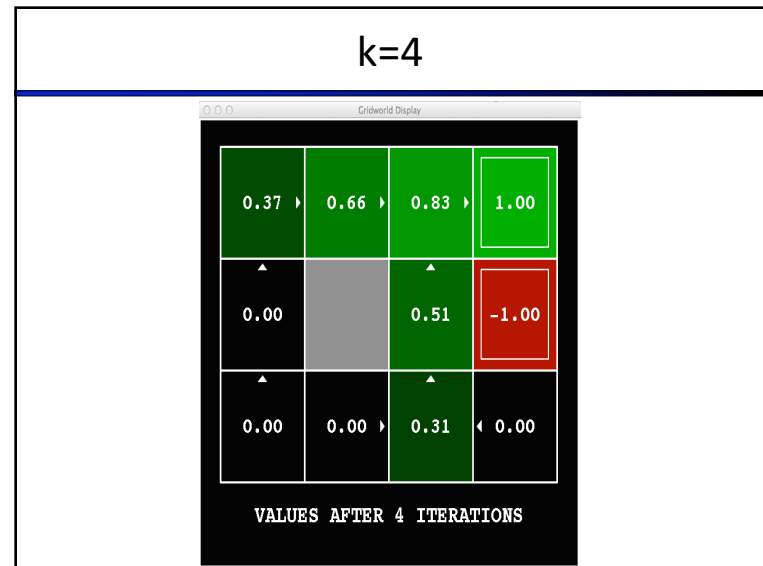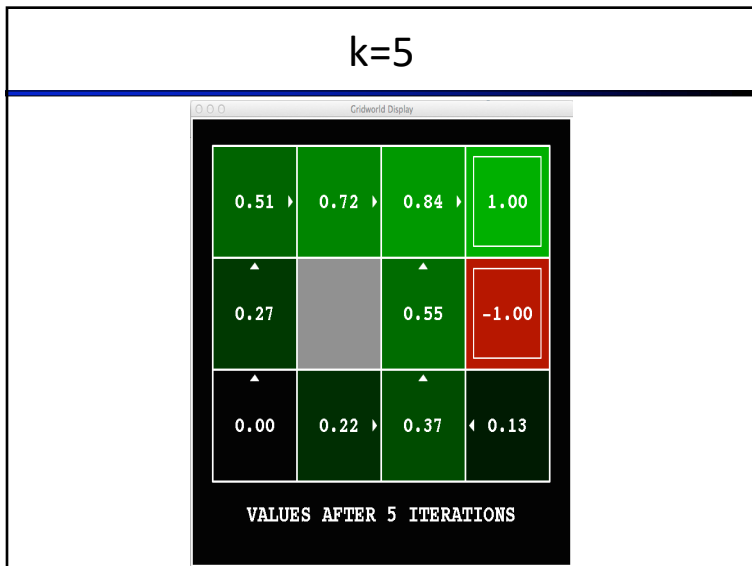
10

## k=1



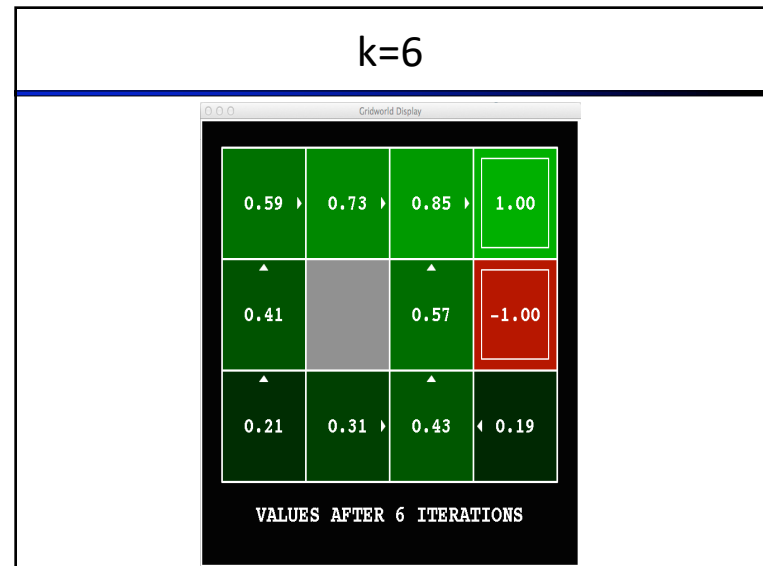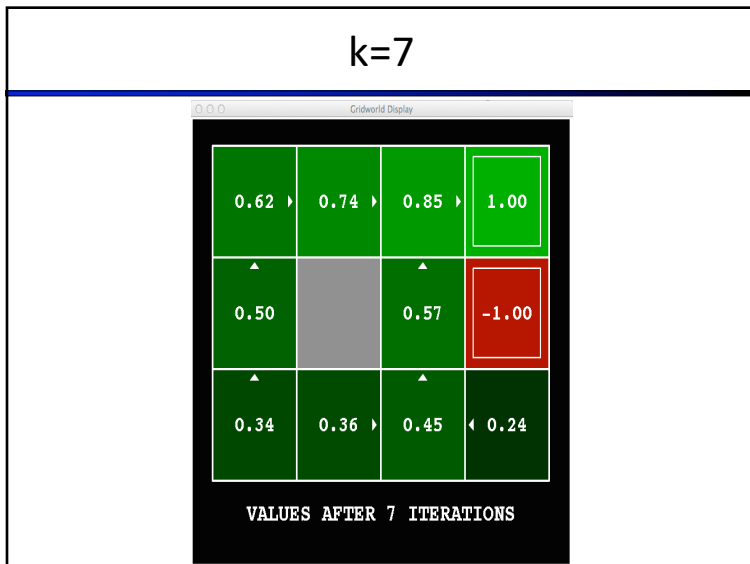Noise = 0.2
Discount = 0.9
Living reward = 0

11

## k=2



12

3

k=3

VALUES AFTER 3 ITERATIONS



k=4

VALUES AFTER 4 ITERATIONS



k=5

VALUES AFTER 5 ITERATIONS



k=6

VALUES AFTER 6 ITERATIONS

13

14

15

16

# k=7



VALUES AFTER 7 ITERATIONS

17

# k=8



VALUES AFTER 8 ITERATIONS

18

# k=9



VALUES AFTER 9 ITERATIONS

19

# k=10



VALUES AFTER 10 ITERATIONS

20

k=11

VALUES AFTER 11 ITERATIONS

21



k=12

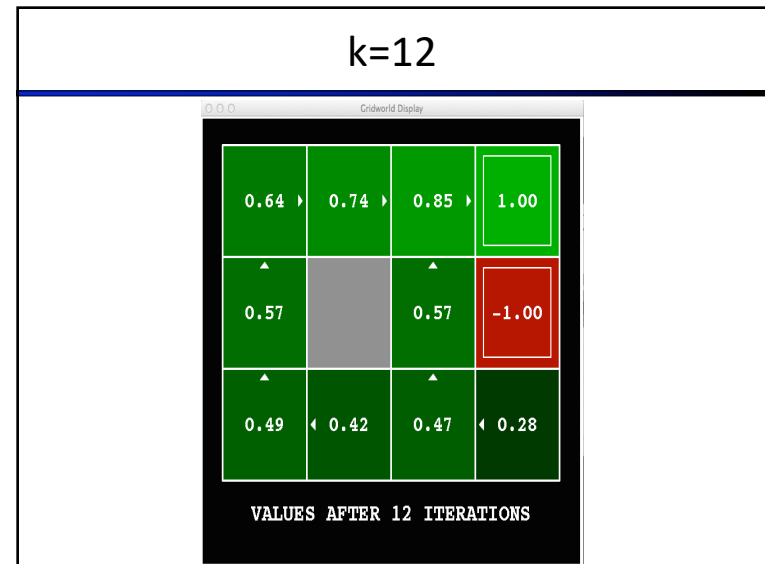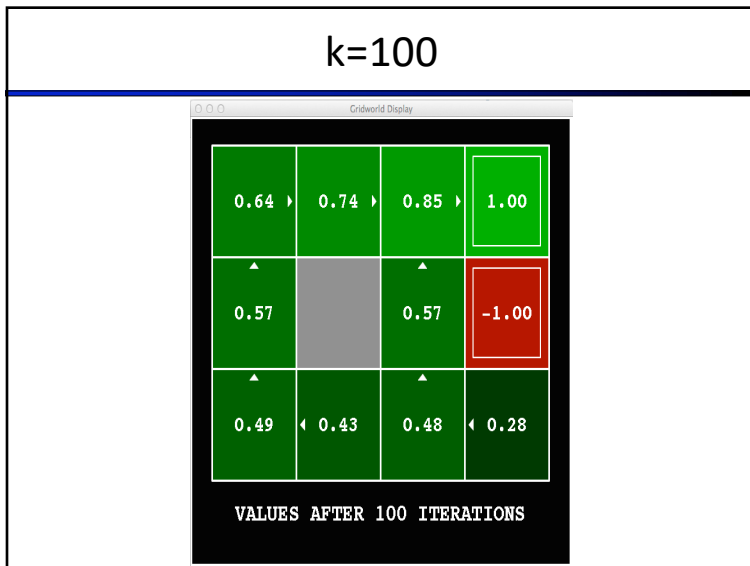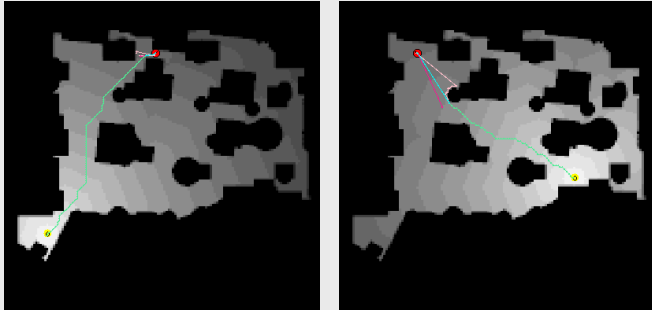VALUES AFTER 12 ITERATIONS

22



k=100

VALUES AFTER 100 ITERATIONS

23

## Value Function and Policy

- Each step takes $O(|A|\ |S|\ |S|)$ time.
- Number of iterations required is polynomial in $|S|$, $|A|$, $1/(1-\text{gamma})$



24

## Value Iteration for Motion Planning
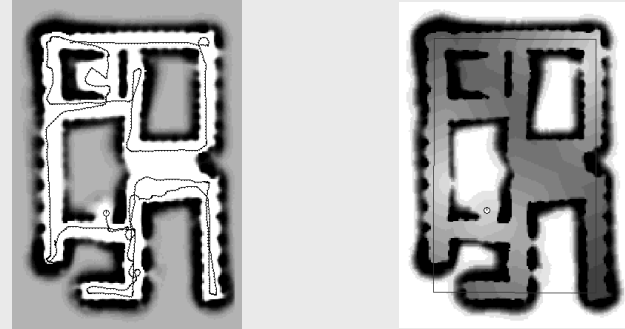**(assumes knowledge of robot's location)**



25

## Frontier-based Exploration

- Every unknown location is a target point.



26

## POMDPs

- In POMDPs we apply the very same idea as in MDPs.
- Since the **state is not observable**, the agent has to **make its decisions based on the belief state** which is a posterior distribution over states.
- For finite horizon problems, the resulting value functions are piecewise linear and convex.
- In each iteration the **number of linear constraints grows exponentially**.
- Full fledged POMDPs have only been applied to very small state spaces with small numbers of possible observations and actions.
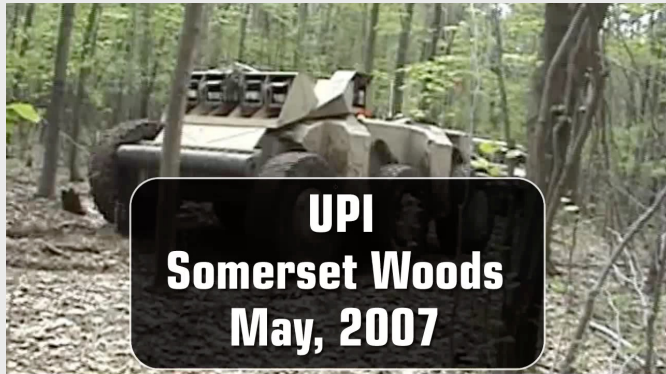- **Approximate solutions are becoming more and more capable**.

27

## CSE 571
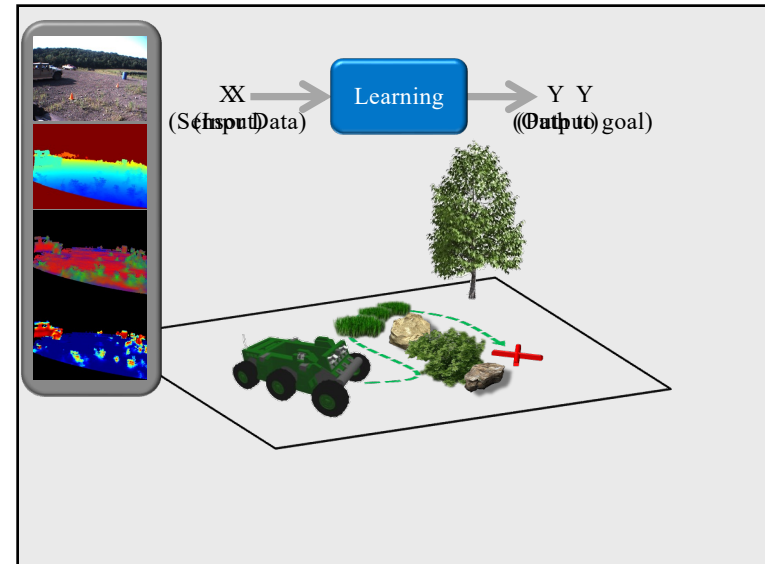## Inverse Optimal Control
## (Inverse Reinforcement Learning)
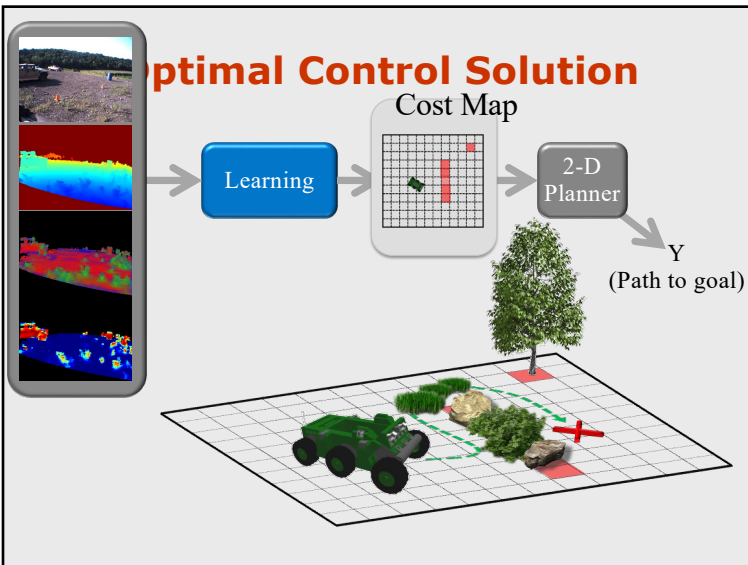
Many slides by Drew Bagnell
Carnegie Mellon University

28

7

**Autonomous Navigation**

UPI
Somerset Woods
May, 2007

29



X X
(Sensor Data) (Input)
→ Learning →
Y Y
(Output) (Path to goal)

30



ptimal Control Solution

Cost Map

Learning → 2-D Planner
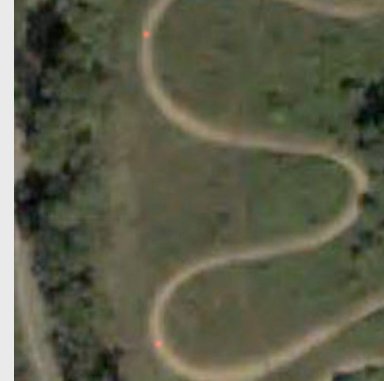
Y
(Path to goal)

31

**Mode 1: Training example**



32

## Mode 1: Training example

## Mode 1: Learned behavior

## Mode 1: Learned behavior
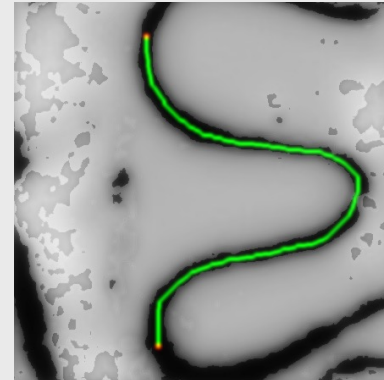
## Mode 1: Learned cost map

## Mode 2: Training example
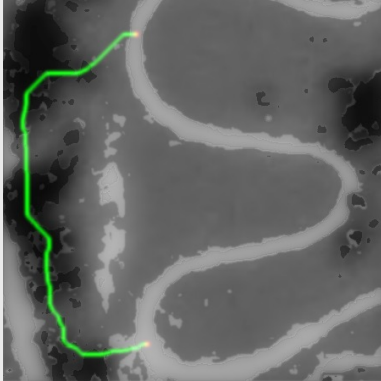


37

## Mode 2: Training example



38

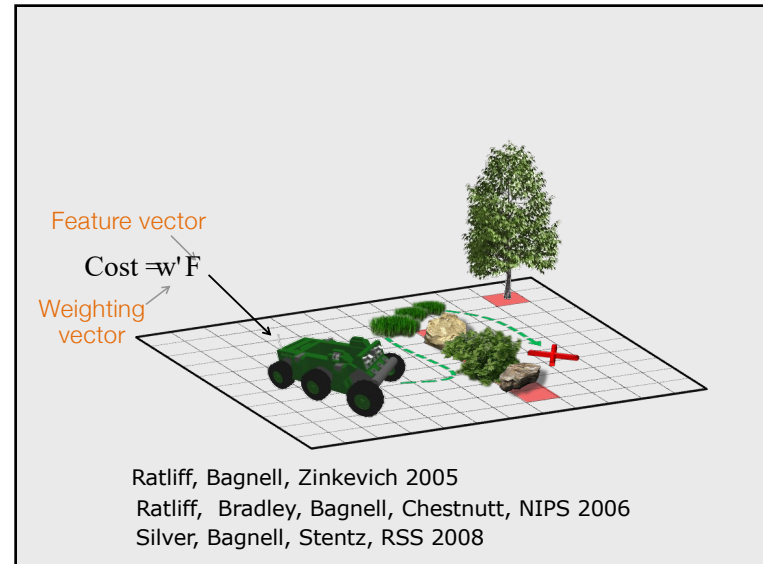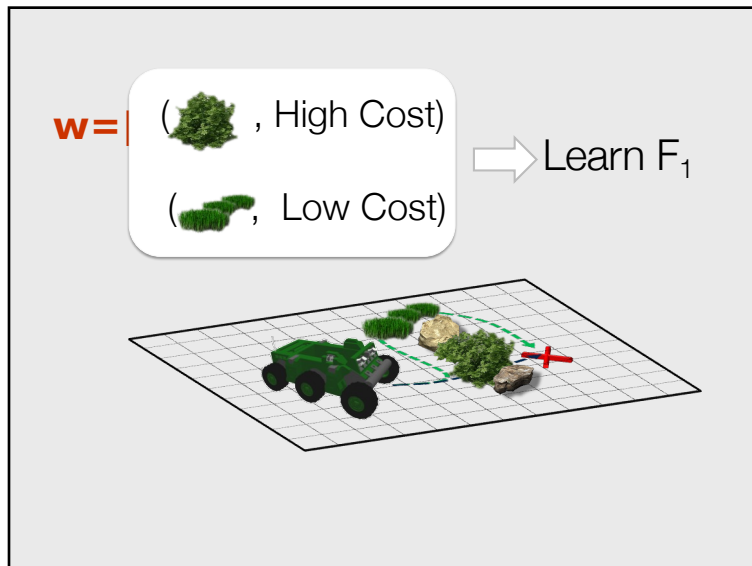## Mode 2: Learned behavior



39

## Mode 2: Learned behavior



40

**Mode 2: Learned cost map**

41

Feature vector

$Cost = w'F$

Weighting vector

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

42

**w=[** ( 🌳 , High Cost)

( 🌿 , Low Cost) ⟹ Learn $F_1$

43

**w=[** ( 🪨 , High Cost)

( 🌿 , Low Cost) ⟹ Learn $F_2$

44

example path

45

---

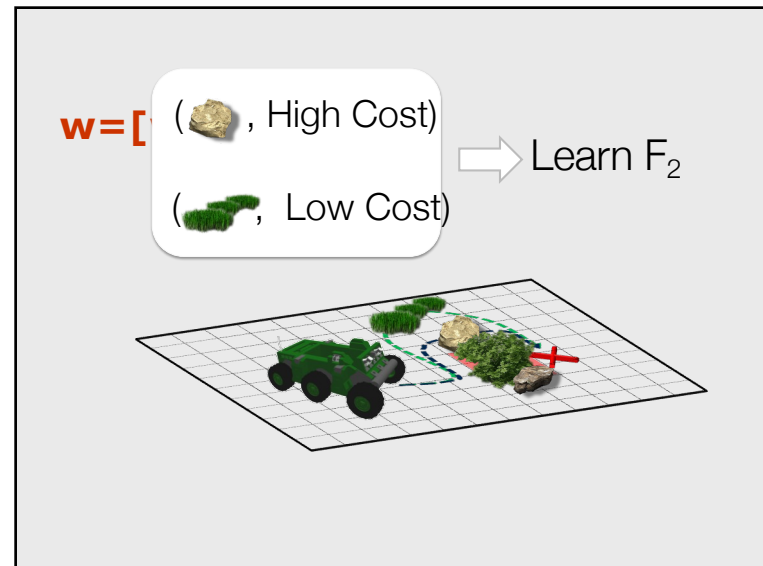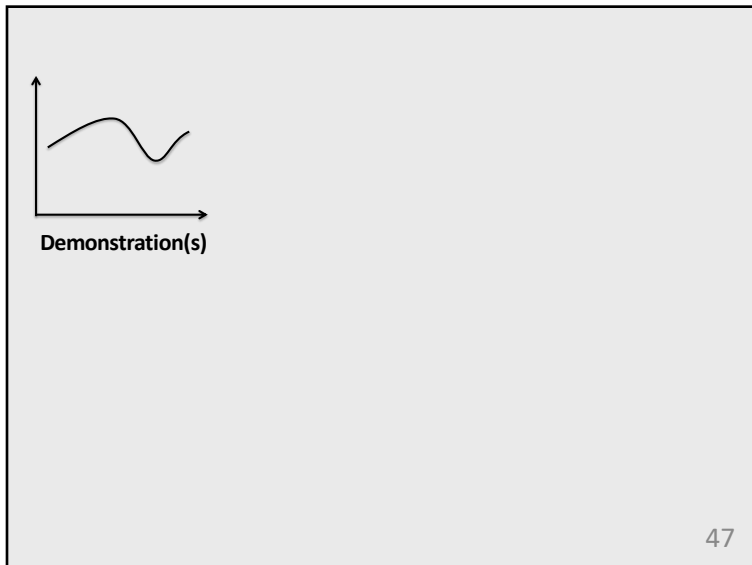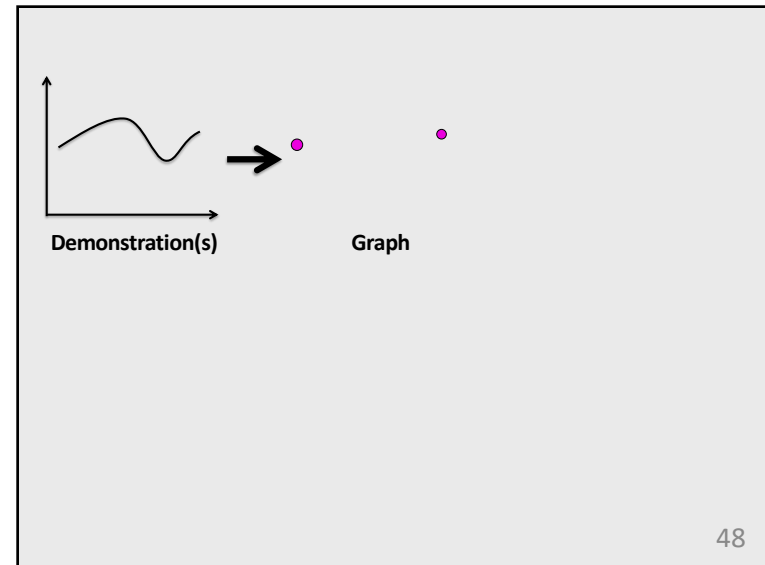**Learning Manipulation Preferences**

- *Input:* Human demonstrations of preferred behavior (e.g., moving a cup of water upright without spilling)

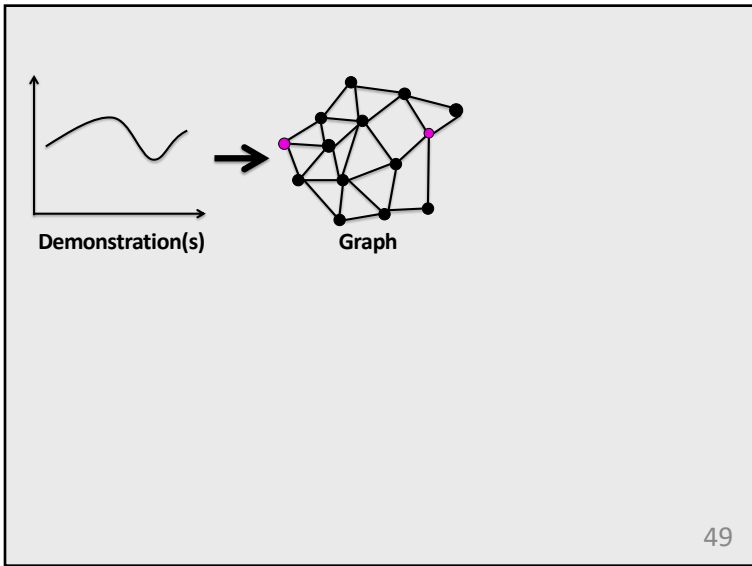- *Output:* Learned cost function that results in trajectories satisfying user preferences



46

46

---



Demonstration(s)

47

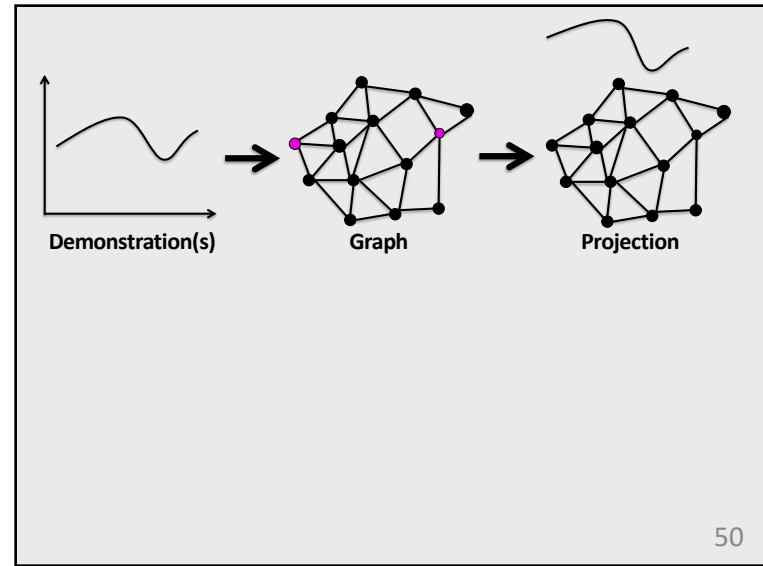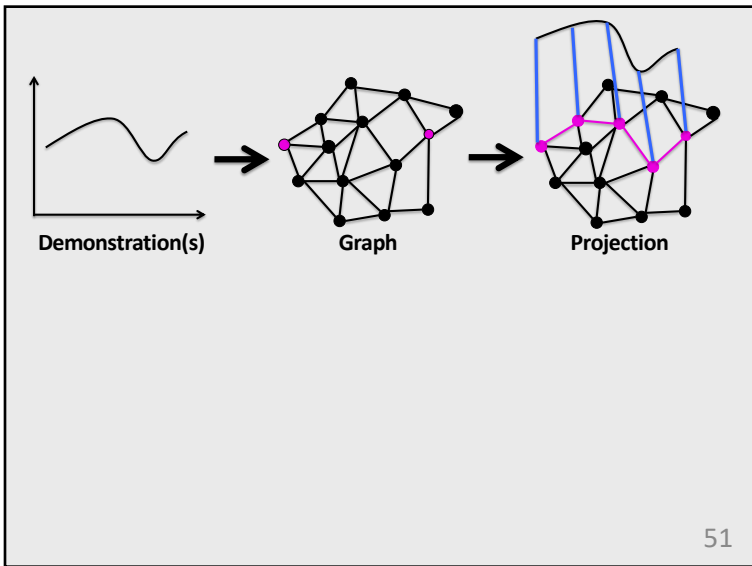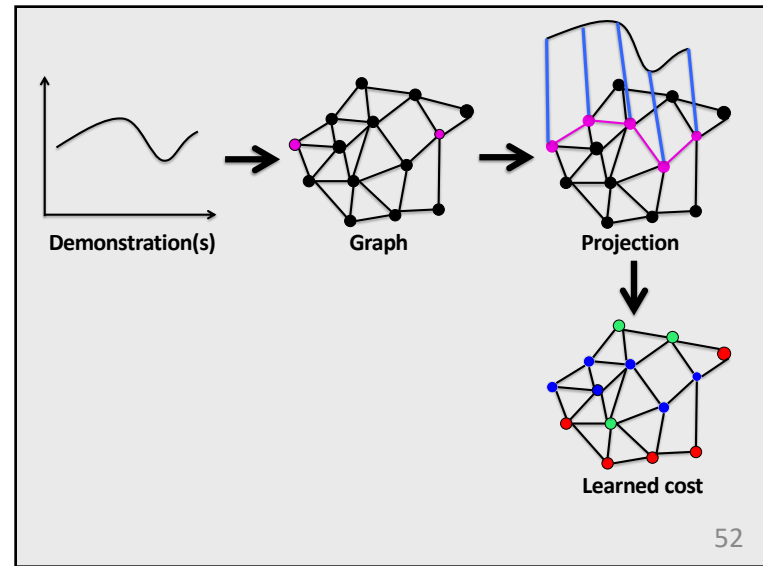47

---



Demonstration(s)          Graph

48

48

---

Demonstration(s)     Graph

49



Demonstration(s)     Graph     Projection

50



Demonstration(s)     Graph     Projection

51



Demonstration(s)     Graph     Projection

Learned cost

52

Demonstration(s) — Graph — Projection — Learned cost — Discrete sampled paths

53



Demonstration(s) — Graph — Projection — Learned cost — Discrete sampled paths — Output trajectories

54



Demonstration(s) — Graph — Projection — Discrete MaxEnt IOC — Learned cost — Discrete sampled paths — Output trajectories

55



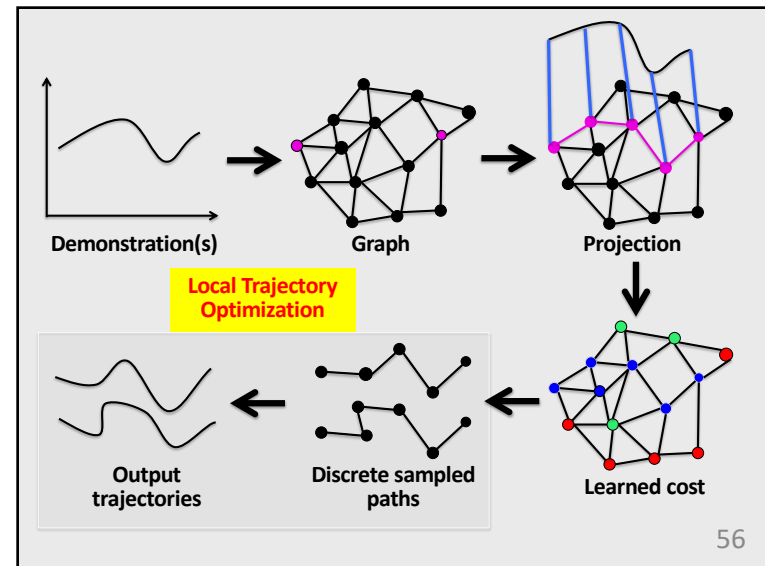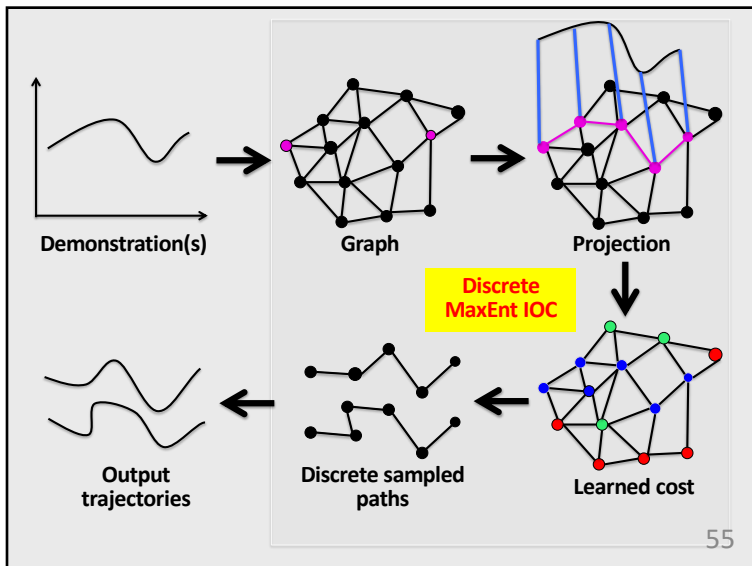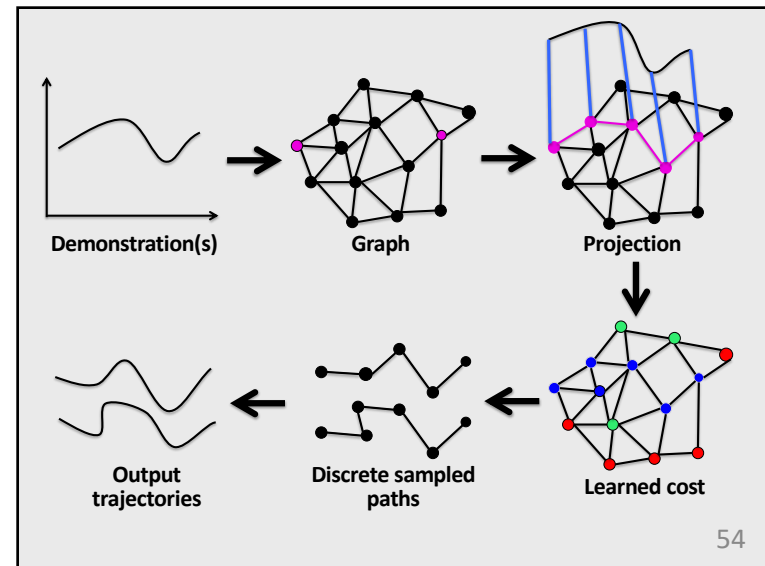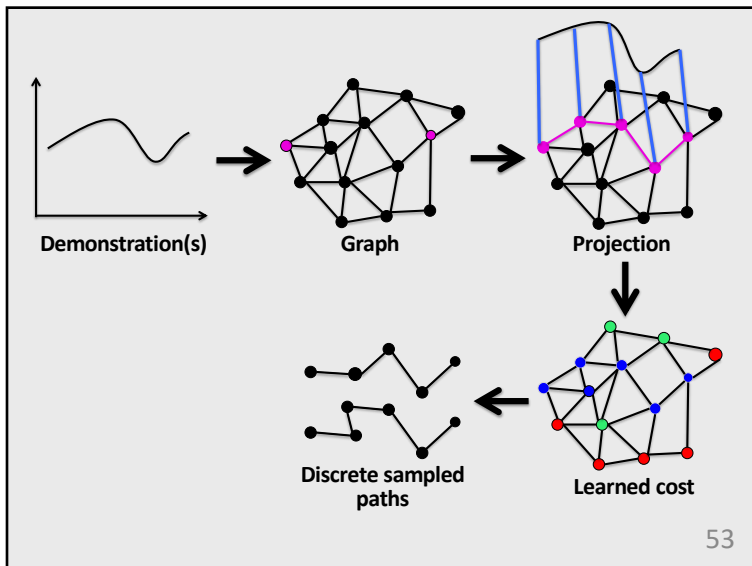Demonstration(s) — Graph — Projection — Local Trajectory Optimization — Learned cost — Discrete sampled paths — Output trajectories

56

## Setup

- *Binary* state-dependent features (~95)
  - Histograms of distances to objects
  - Histograms of end-effector orientation
  - Object specific features (electronic vs non-electronic)
  - Approach direction w.r.t goal

- **Task**
  - Hold cup upright while not moving above electronics

## Laptop task: Demonstration
**( Not part of training set)**

## Laptop task: LTO + Smooth random path

## Readings

- Max-Ent IRL (Ziebart, Bagnell): http://www.cs.cmu.edu/~bziebart/
- CIOC (Levine) http://graphics.stanford.edu/projects/cioc/cioc.pdf
- Manipulation (Byravan/Fox): https://rse-lab.cs.washington.edu/papers/graph-based-IOC-ijcai-2015.pdf
- Imitation learning (Ermon): https://cs.stanford.edu/~ermon/
- Human/manipulation (Dragan): https://people.eecs.berkeley.edu/~anca/research.html