

# CSE-573 Artificial Intelligence

## **Partially-Observable MDPS (POMDPs)**

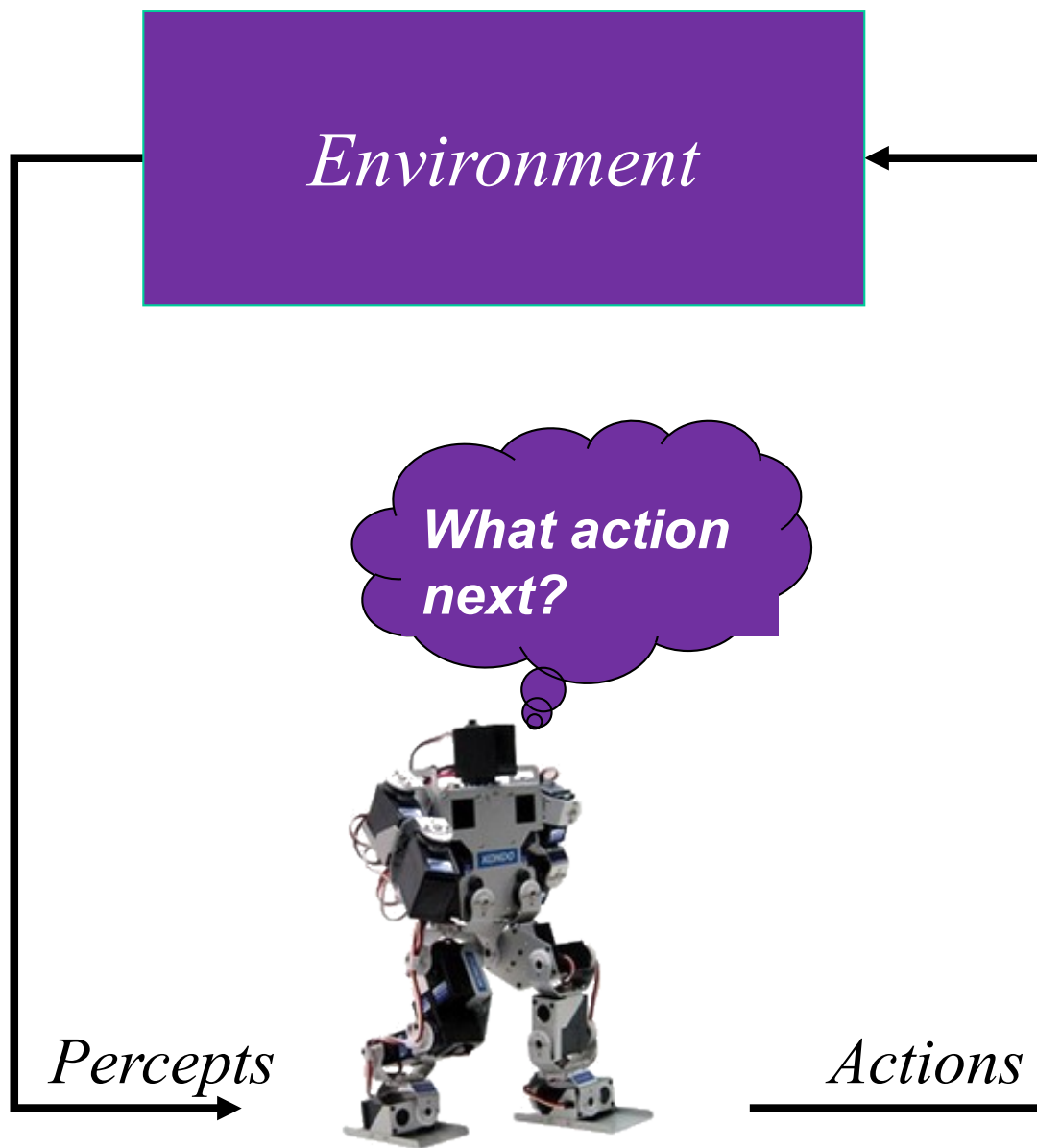
# Todo

- Key slides don't have Y axis labeled – NOT value

# Classical

**Fully  
Observable**

**Perfect**



**Deterministic**

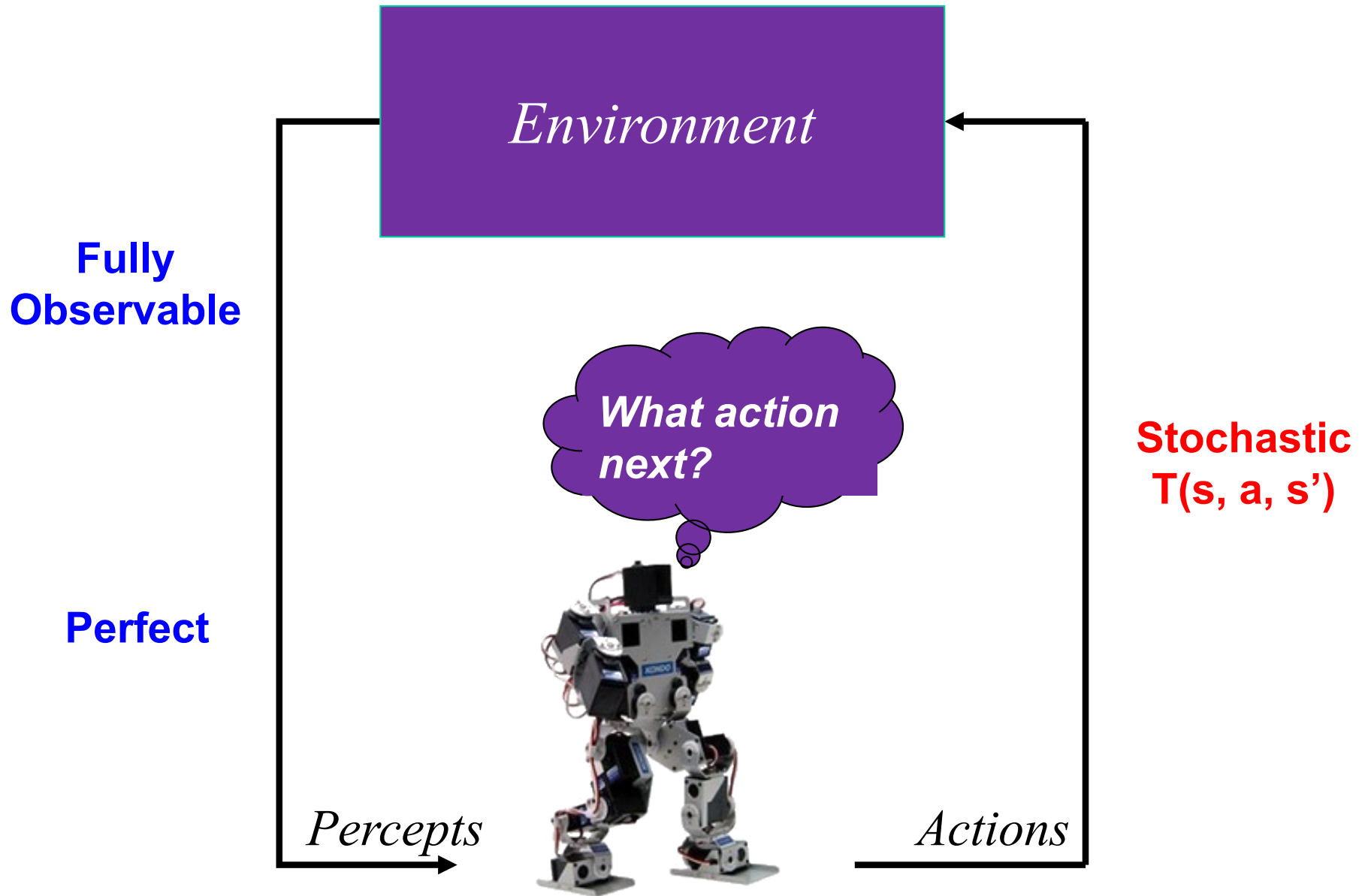
*Actions*

*Percepts*

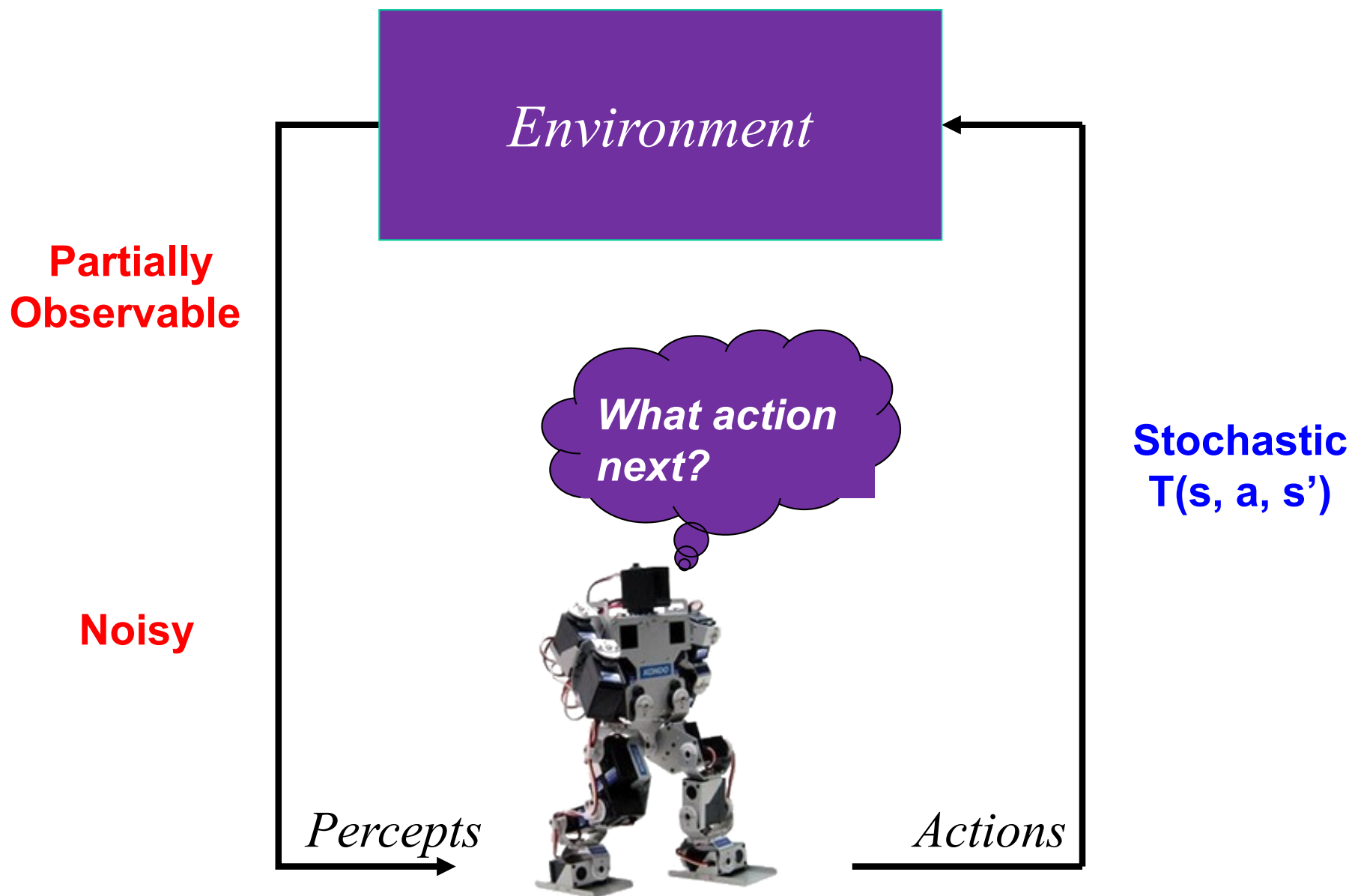
*Environment*

*What action  
next?*

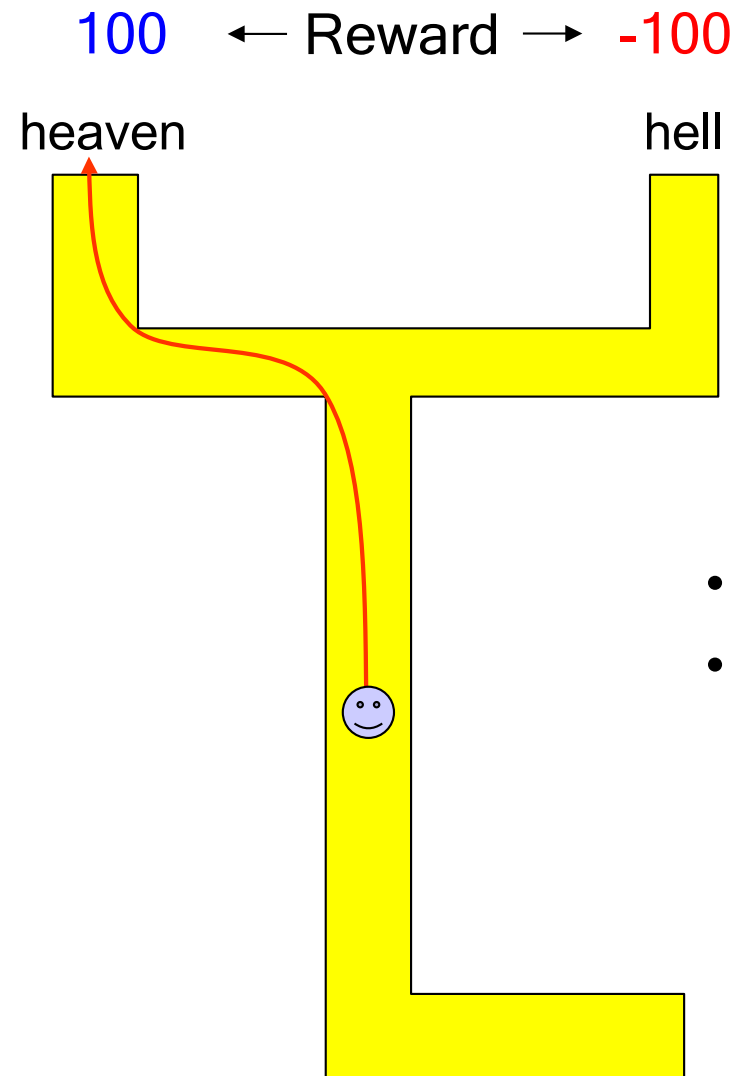
# Stochastic (MDP)



# Partially-Observable Stochastic (POMDP)



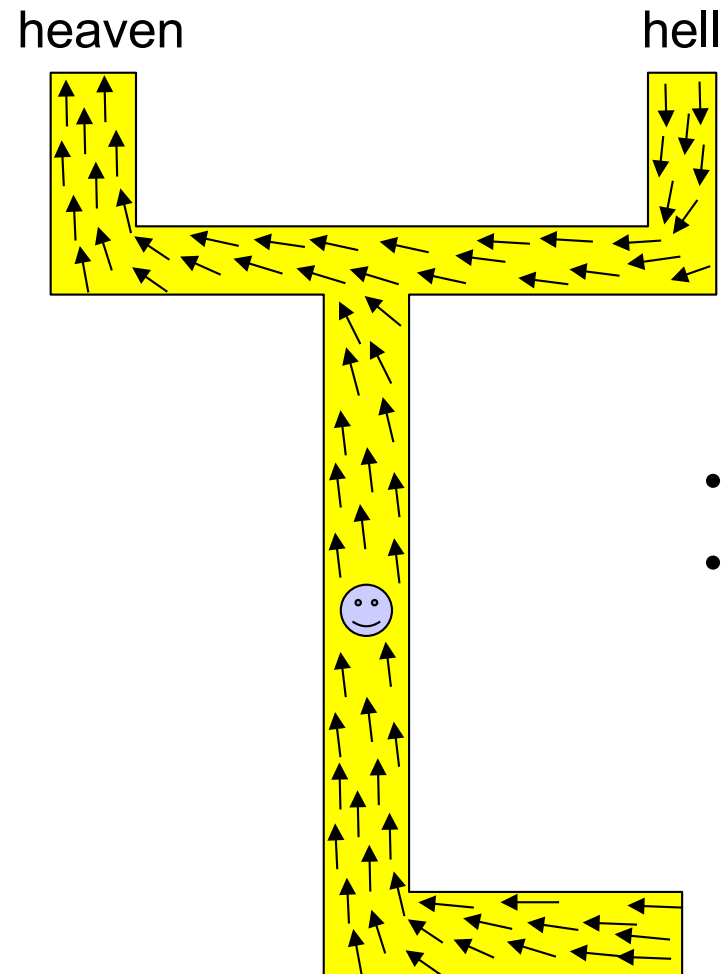
# Classical Planning



- Sequential Plan

- World deterministic
- State observable

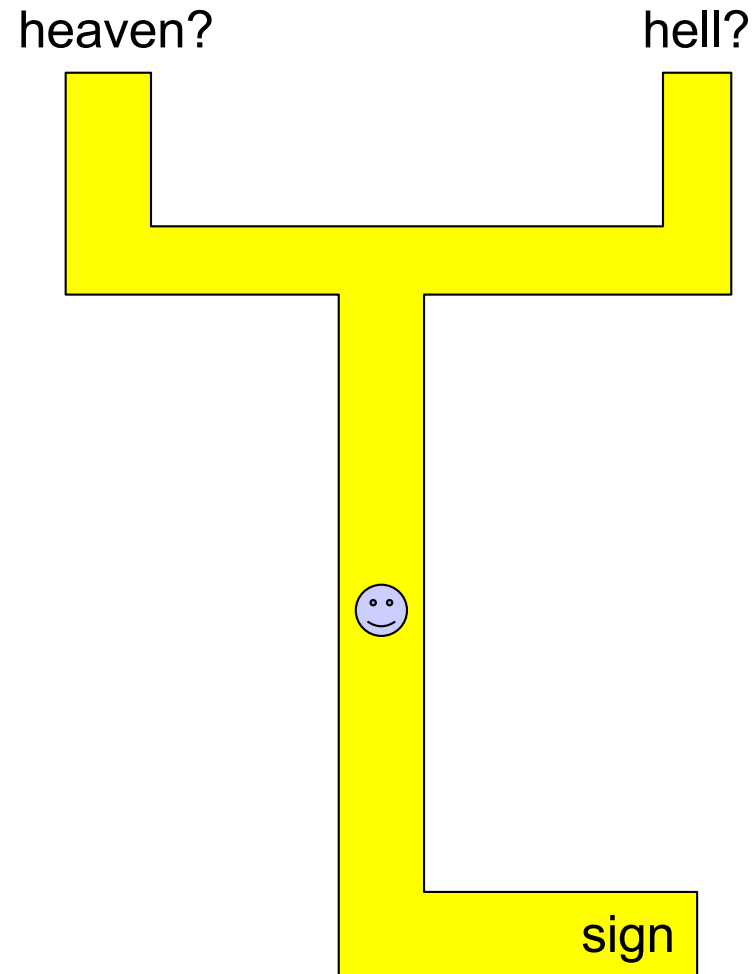
# MDP-Style Planning



- Policy

- World stochastic
- State observable

# Stochastic, *Partially* Observable





# Markov Decision Process (MDP)

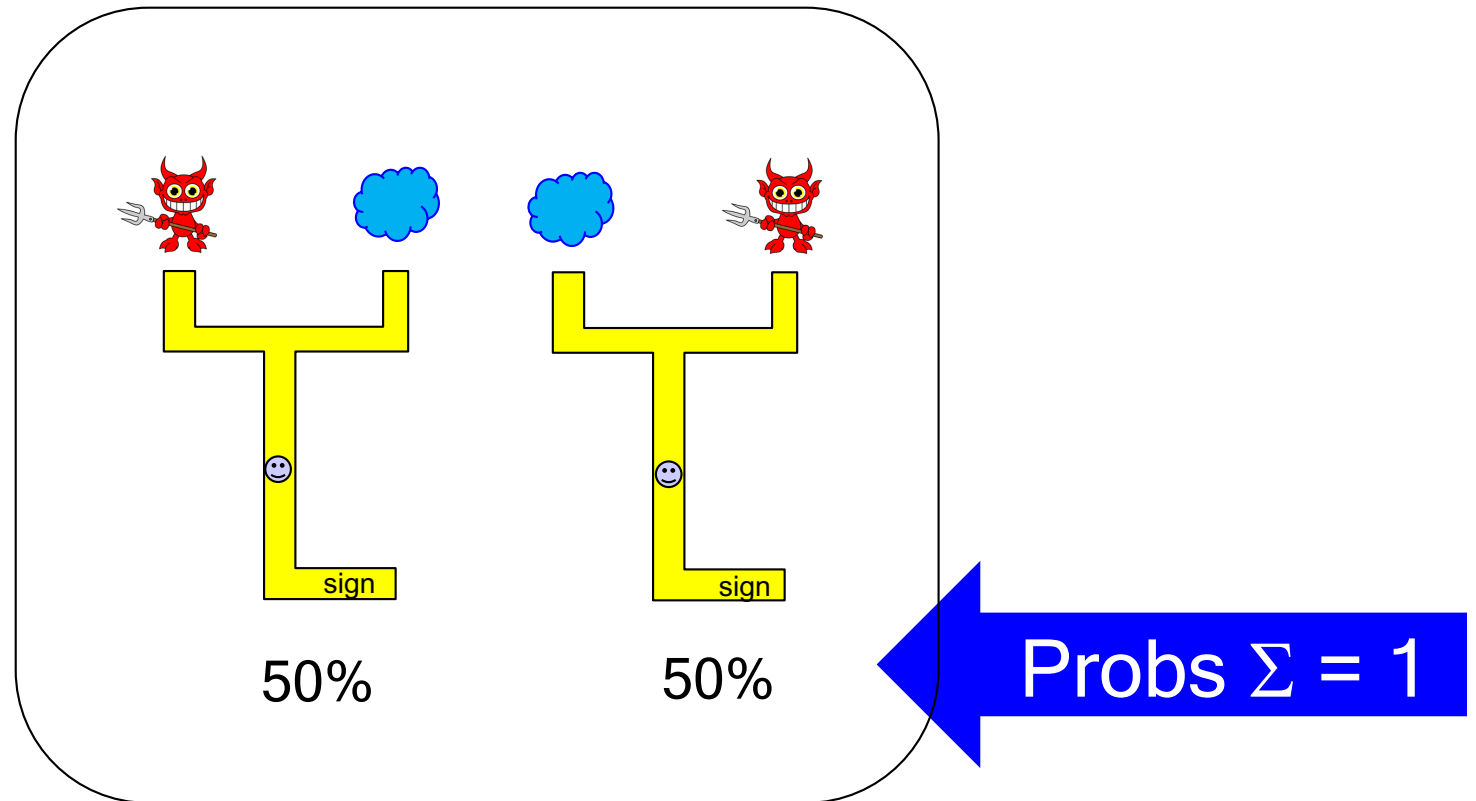
- **S**: set of states
- **A**: set of actions
- $\Pr(s' | s, a)$ : transition model
- $\mathbf{R}(s, a, s')$ : reward model
- $\gamma$ : discount factor
- $s_0$ : start state

# Partially-Observable MDP

- **S**: set of states
- **A**: set of actions
- $\Pr(s' | s, a)$ : transition model
- $\mathbf{R}(s, a, s')$ : reward model
- $\gamma$ : discount factor
- $s_0$ : start state
- **E** set of possible evidence (observations)
- $\Pr(e | s)$

# Belief State

- State of agent's mind
- Not just of world

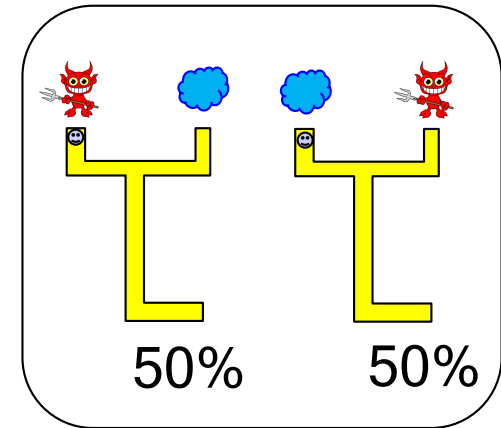
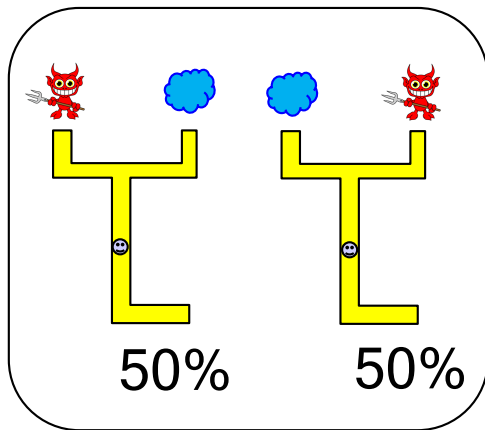


Note: POMDP

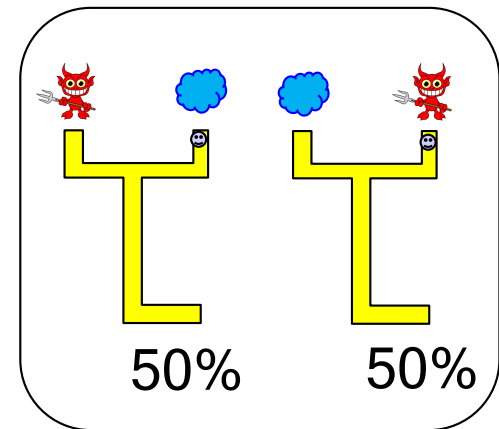
# Planning in Belief Space

For now, assume movement is deterministic

And *NO* observations possible



Exp. Reward: 0



Exp. Reward: 0

# Partially-Observable MDP

- **S**: set of states
- **A**: set of actions
- $\Pr(s' | s, a)$ : transition model
- $\mathbf{R}(s, a, s')$ : reward model
- $\gamma$ : discount factor
- $s_0$ : start state
- **E** set of possible evidence (aka observations, measurements)
- $\Pr(e | s)$

# Evidence Model

e/w = location of devil

b/m/ul/ur = location of agent

■  $S = \{s_{wb}, s_{eb}, s_{wm}, s_{em}, s_{wul}, s_{eul}, s_{wur}, s_{eur}\}$

■  $E = \{\text{heat}\}$

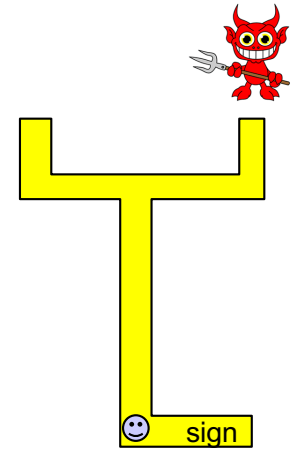
■  $\Pr(e|s)$ :

$$\Pr(\text{heat} | s_{eb}) = 1.0$$

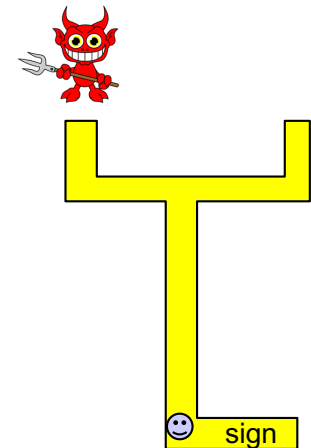
$$\Pr(\text{heat} | s_{wb}) = 0.2$$

$$\Pr(\text{heat} | s_{\text{other}}) = 0.0$$

$s_{eb}$



$s_{wb}$

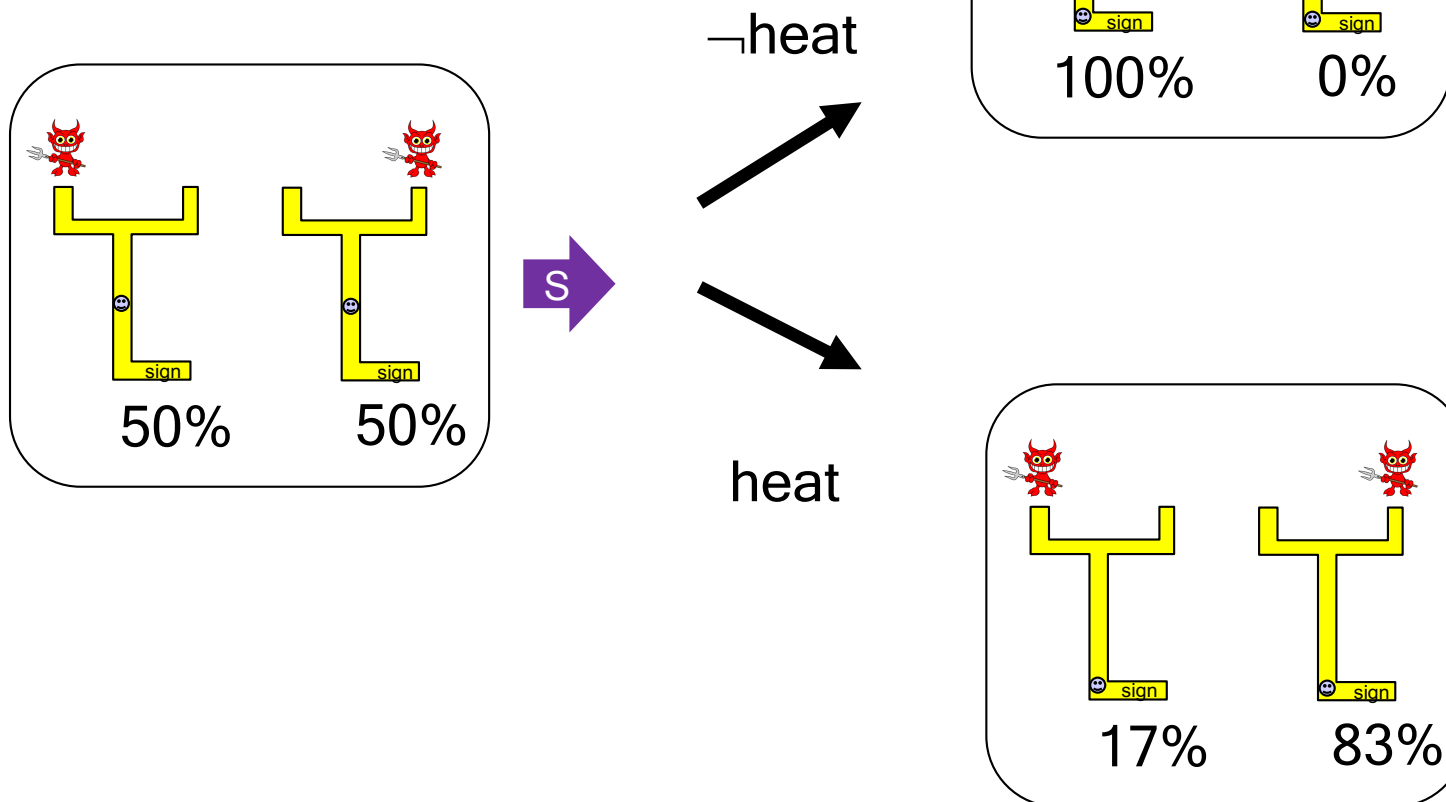


# Updating beliefs given evidence

$$\Pr(\text{heat} \mid s_{\text{eb}}) = 1.0$$
$$\Pr(\text{heat} \mid s_{\text{wb}}) = 0.2$$

Use Bayes rule:

$$P(s \mid e) = P(e \mid s)P(s) / P(e)$$



# Objective of a Fully Observable MDP

- Find a policy

$$\pi: \mathbf{S} \rightarrow \mathbf{A}$$

- which maximizes expected discounted reward
  - given an infinite horizon
  - assuming full observability



# Objective of a POMDP

- Find a policy

$$\pi: \text{BeliefStates}(\mathbf{S}) \rightarrow \mathbf{A}$$

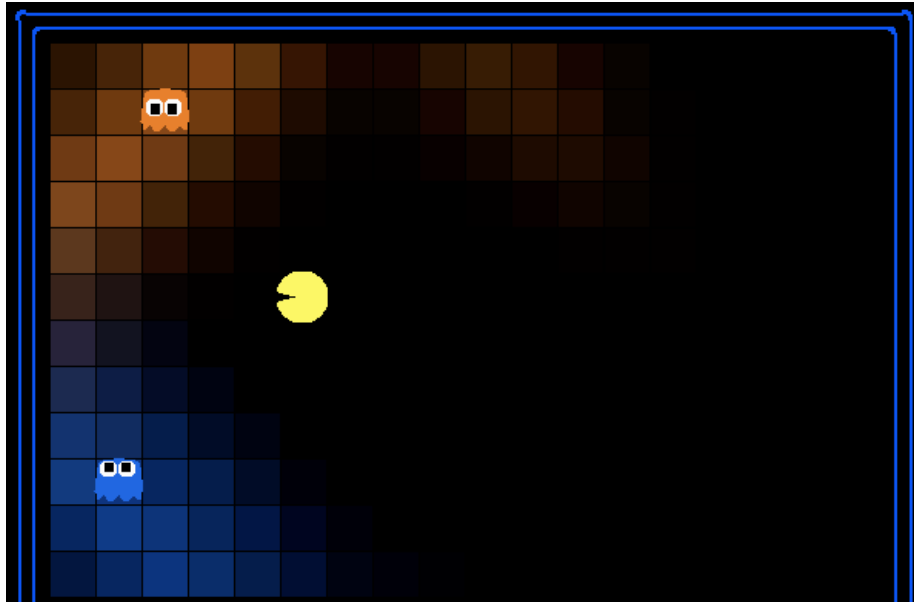
A belief state is a *probability distribution* over states

- which maximizes expected discounted reward

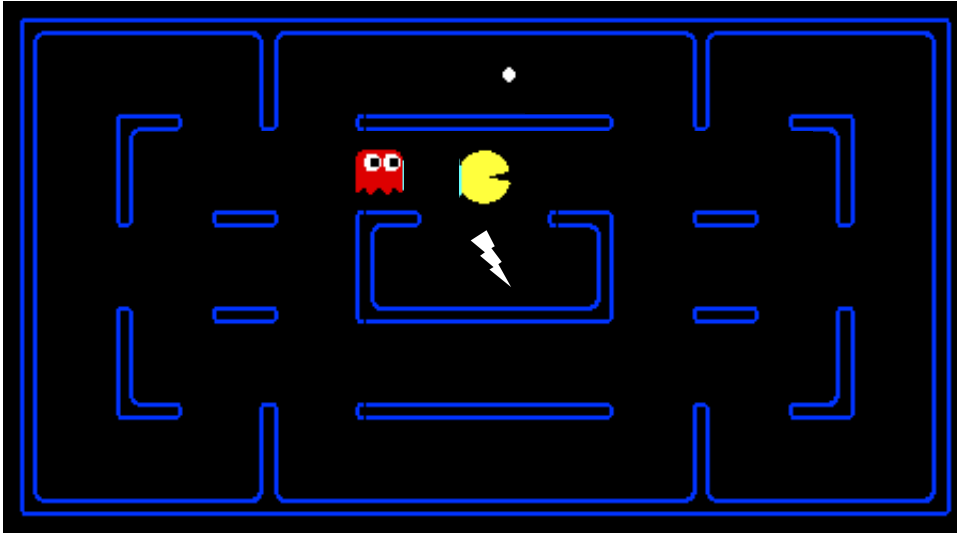
- given an infinite horizon
- assuming *partial* & *noisy* observability

# Planning in last HW

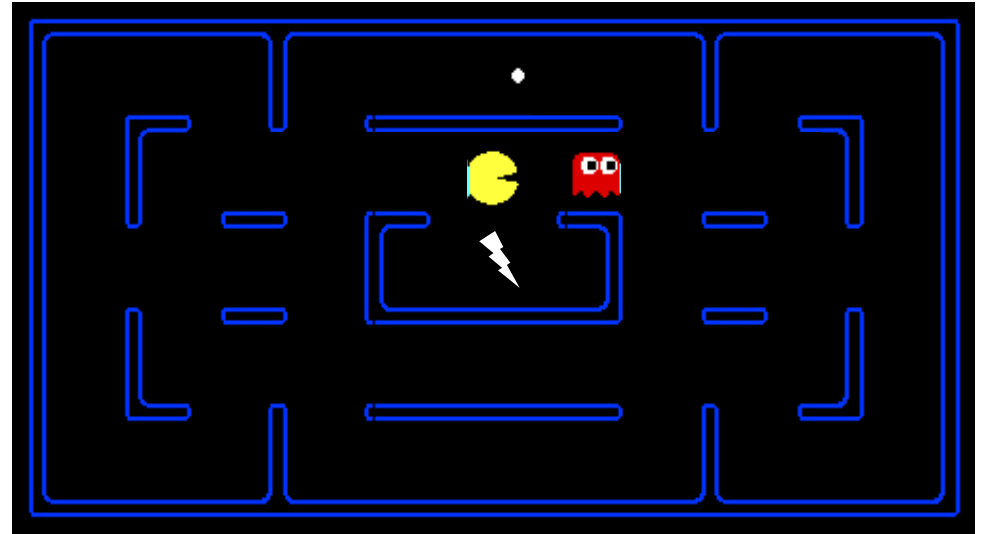
- Map Estimate
- Now “know” state
- Solve MDP



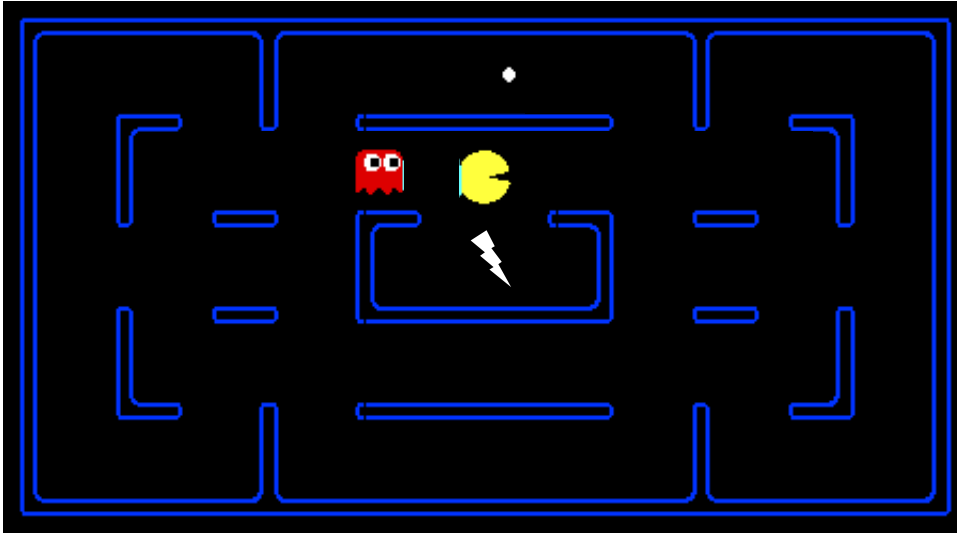
Best plan to eat final food?



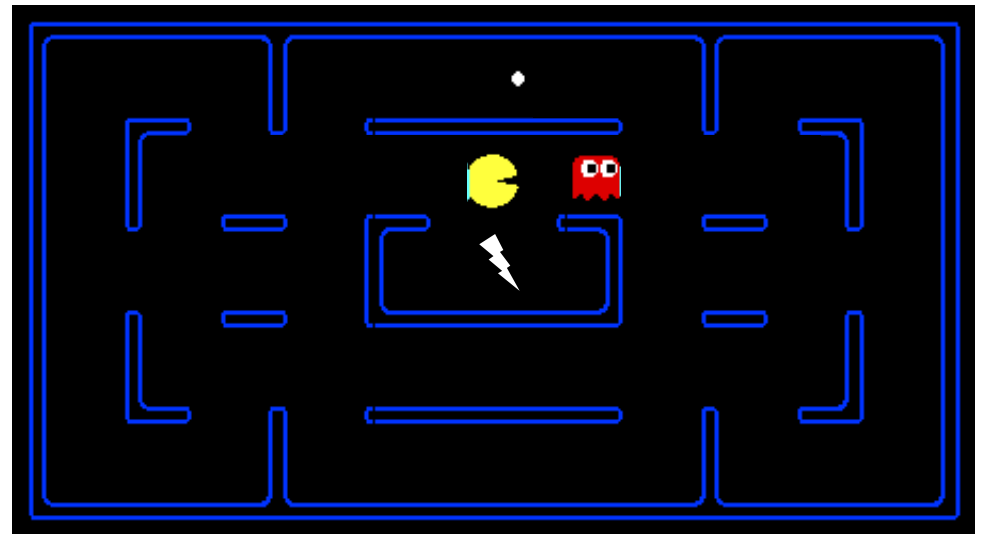
Best plan to eat final food?



# Problem with Planning from MAP Estimate



49%



51%

- Best action for belief state over  $k$  worlds may not be the best action in *any one* of those worlds

# 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*.

$\pi : \text{beliefs} \rightarrow \text{actions}$

- Let  $b$  be the belief of the agent about the state under consideration.

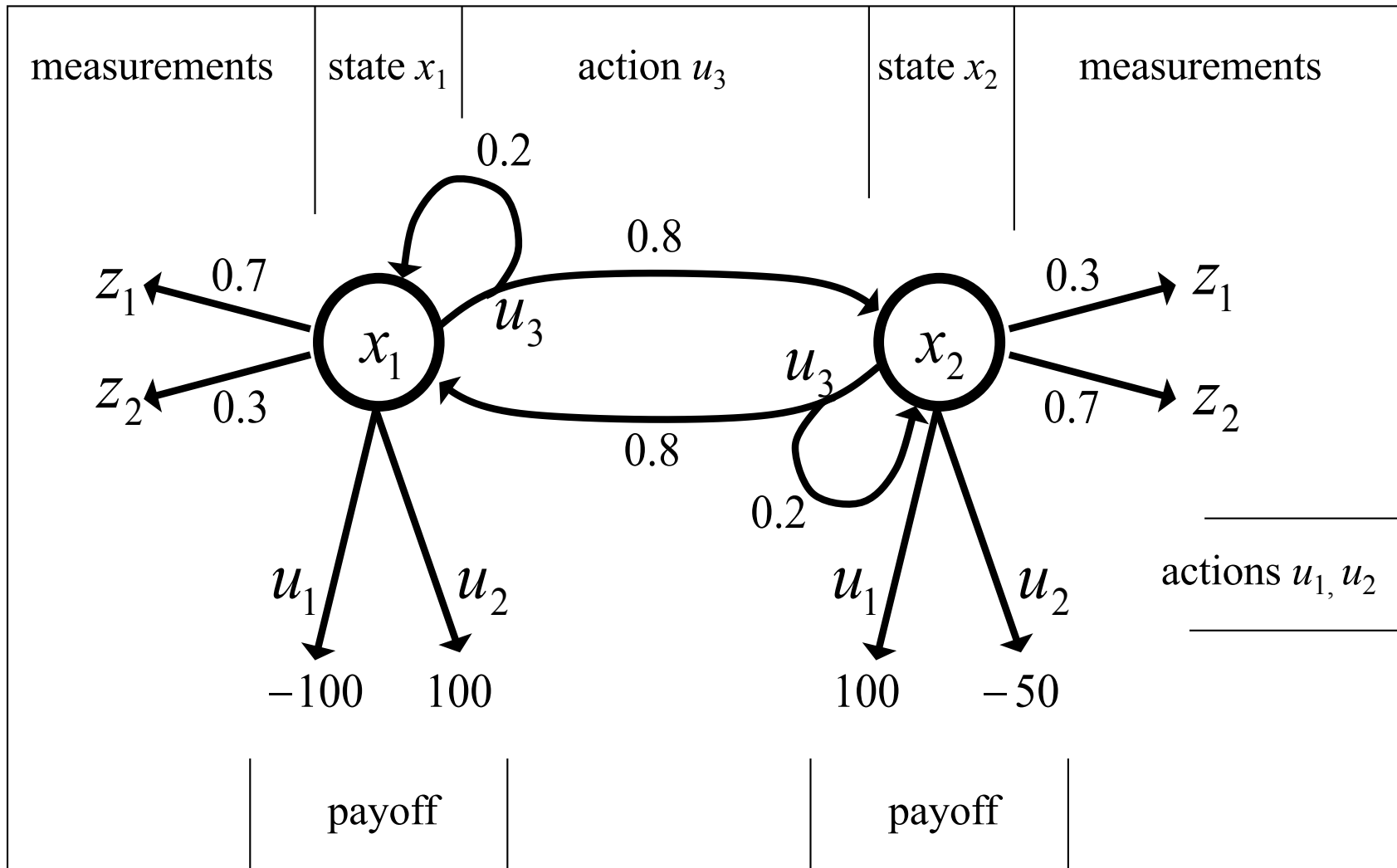
- POMDPs compute a *value function over belief space*:

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

# Problems

- Each belief is a probability distribution, thus, each value in a **POMDP is a function of an entire probability distribution.**
- **This is challenging, since probability distributions are continuous.**
  - How many belief states are there?
  - How many policies are there?
- For **finite worlds** with finite state, action, and evidence spaces and finite horizons, however, we can **effectively represent the value functions by piecewise linear functions.**

# An Illustrative Example





# The Parameters of the Example

- The actions  $u_1$  and  $u_2$  are terminal actions.
- The action  $u_3$  is a sensing action that potentially leads to a state transition.
- The horizon is finite and  $\gamma=1$ .

$$r(x_1, u_1) = -100$$

$$r(x_2, u_1) = +100$$

$$r(x_1, u_2) = +100$$

$$r(x_2, u_2) = -50 \leftarrow$$

$$r(x_1, u_3) = -1$$

$$r(x_2, u_3) = -1$$

$$p(x'_1|x_1, u_3) = 0.2$$

$$p(x'_2|x_1, u_3) = 0.8$$

$$p(x'_1|x_2, u_3) = 0.8$$

$$p(z'_2|x_2, u_3) = 0.2$$

$$p(z_1|x_1) = 0.7$$

$$p(z_2|x_1) = 0.3$$

$$p(z_1|x_2) = 0.3$$

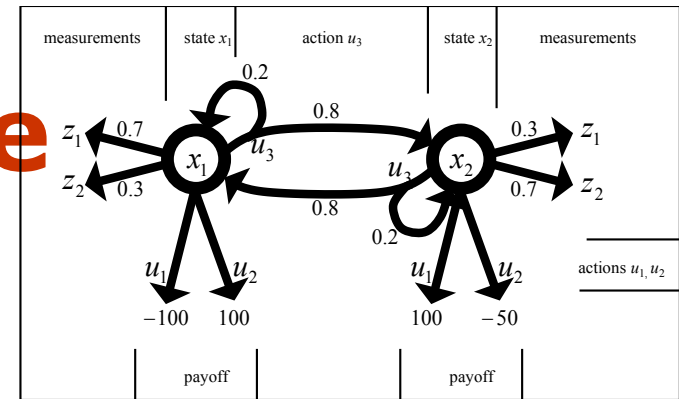
$$p(z_2|x_2) = 0.7$$

# Payoff in POMDPs

- In MDPs, the payoff (or return) depended on the state of the system.
- In POMDPs, however, the true state is not exactly known.
- Therefore, we compute the **expected payoff** by **integrating over all states**:

$$\begin{aligned} r(b, u) &= E_x[r(x, u)] \\ &= \int r(x, u)p(x) dx \\ &= p_1 r(x_1, u) + p_2 r(x_2, u) \end{aligned}$$

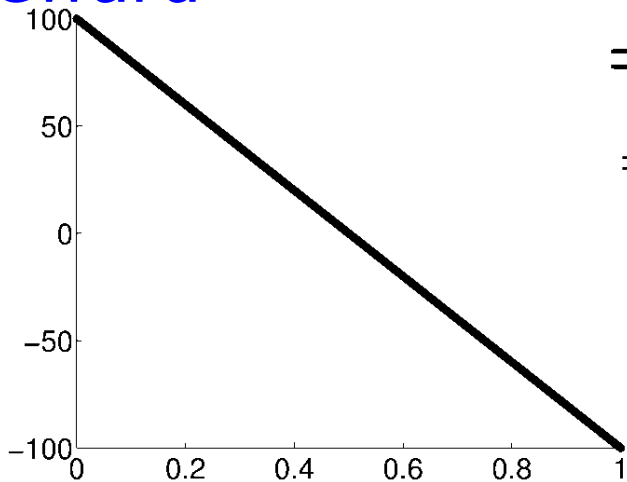
# Payoffs in Our Example



- If we are totally certain that we are in state  $x_1$  and execute action  $u_1$ , we receive a reward of -100
- If, on the other hand, we definitely know that we are in  $x_2$  and execute  $u_1$ , the reward is +100.
- In between it is the linear combination of the extreme values weighted by the probabilities

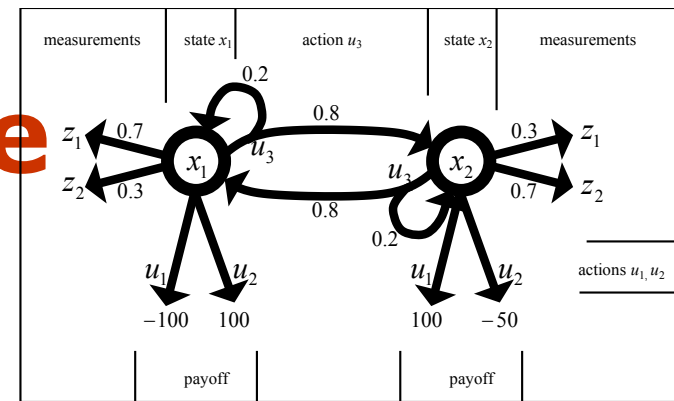
Reward

$$\begin{aligned}
 r(b, u_1) &= -100 p_1 + 100 p_2 \\
 &= -100 p_1 + 100 (1 - p_1) \\
 &= 100 - 200 p_1
 \end{aligned}$$



$P_1 = P(\text{state}=x_1)$

# Payoffs in Our Example



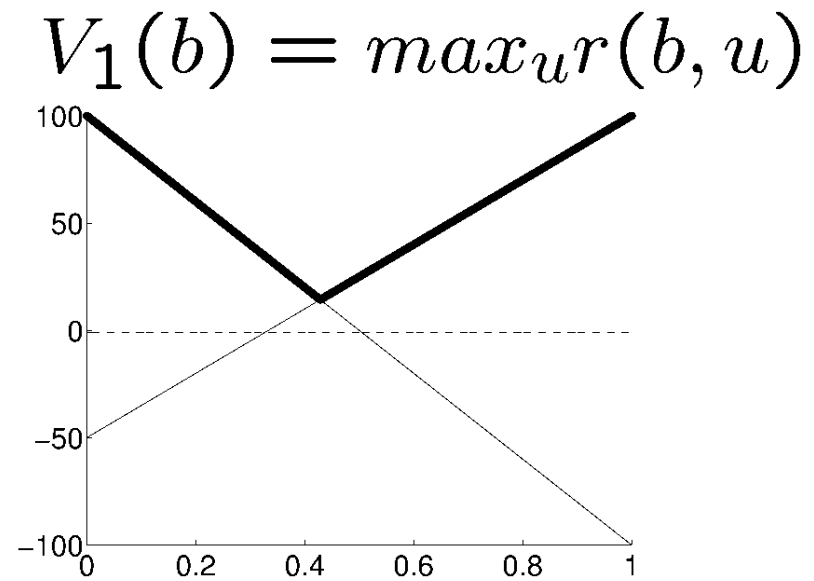
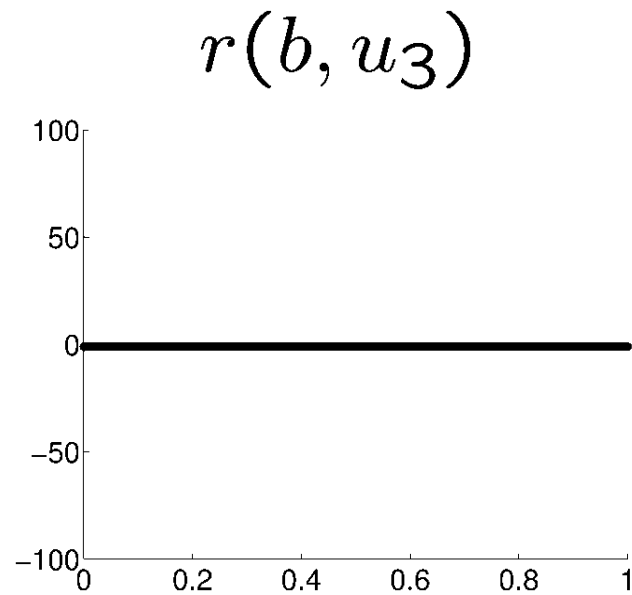
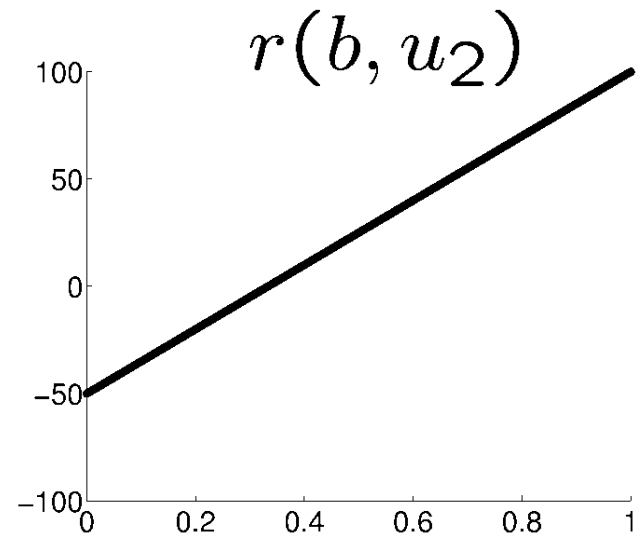
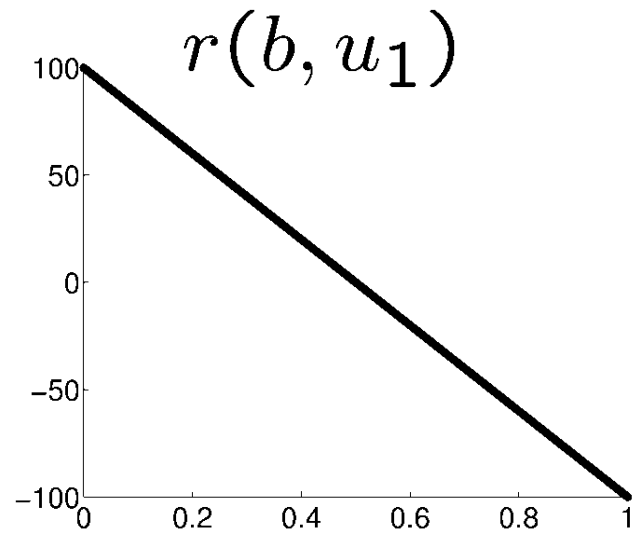
- If we are totally certain that we are in state  $x_1$  and execute action  $u_1$ , we receive a reward of -100
- If, on the other hand, we definitely know that we are in  $x_2$  and execute  $u_1$ , the reward is +100.
- In between it is the linear combination of the extreme values weighted by the probabilities

$$\begin{aligned}
 r(b, u_1) &= -100 p_1 + 100 p_2 \\
 &= -100 p_1 + 100 (1 - p_1) \\
 &= 100 - 200 p_1
 \end{aligned}$$

$$\begin{aligned}
 r(b, u_2) &= 100 p_1 - 50 (1 - p_1) \\
 &= 150 p_1 - 50
 \end{aligned}$$

$$r(b, u_3) = -1$$

# Payoffs in Our Example (2)

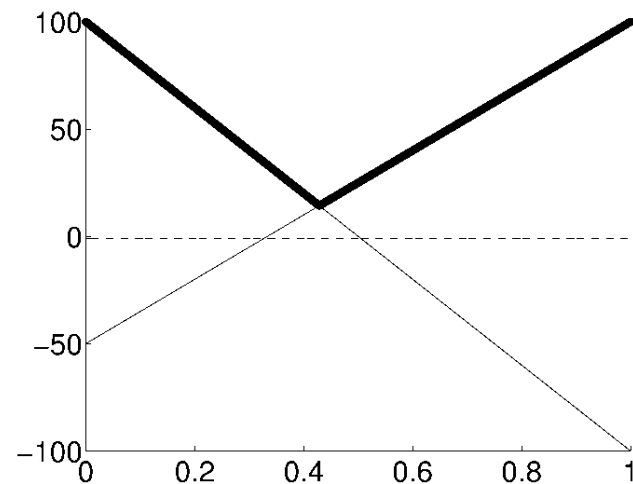


# The Resulting Policy for $T=1$

- Given a finite POMDP with time horizon = 1
- Use  $V_1(b)$  to determine the optimal policy.

$$\pi_1(b) = \begin{cases} u_1 & \text{if } p_1 \leq \frac{3}{7} = 0.429 \\ u_2 & \text{if } p_1 > \frac{3}{7} \end{cases}$$

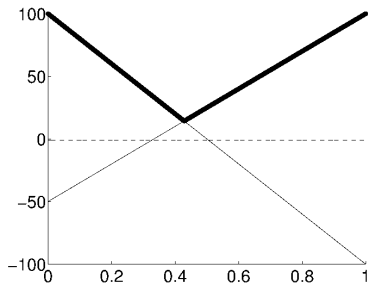
- Corresponding value:



# Piecewise Linearity, Convexity

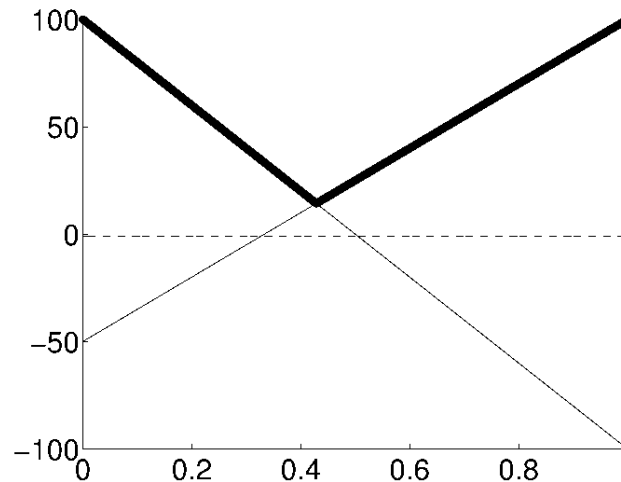
- The resulting value function  $V_1(b)$  is the maximum of the three functions at each point

$$\begin{aligned} V_1(b) &= \max_u r(b, u) \\ &= \max \left\{ \begin{array}{l} -100 p_1 + 100 (1 - p_1) \\ 100 p_1 - 50 (1 - p_1) \\ 0 \end{array} \right\} \end{aligned}$$



- I.e., it's piecewise linear and convex.

# Pruning



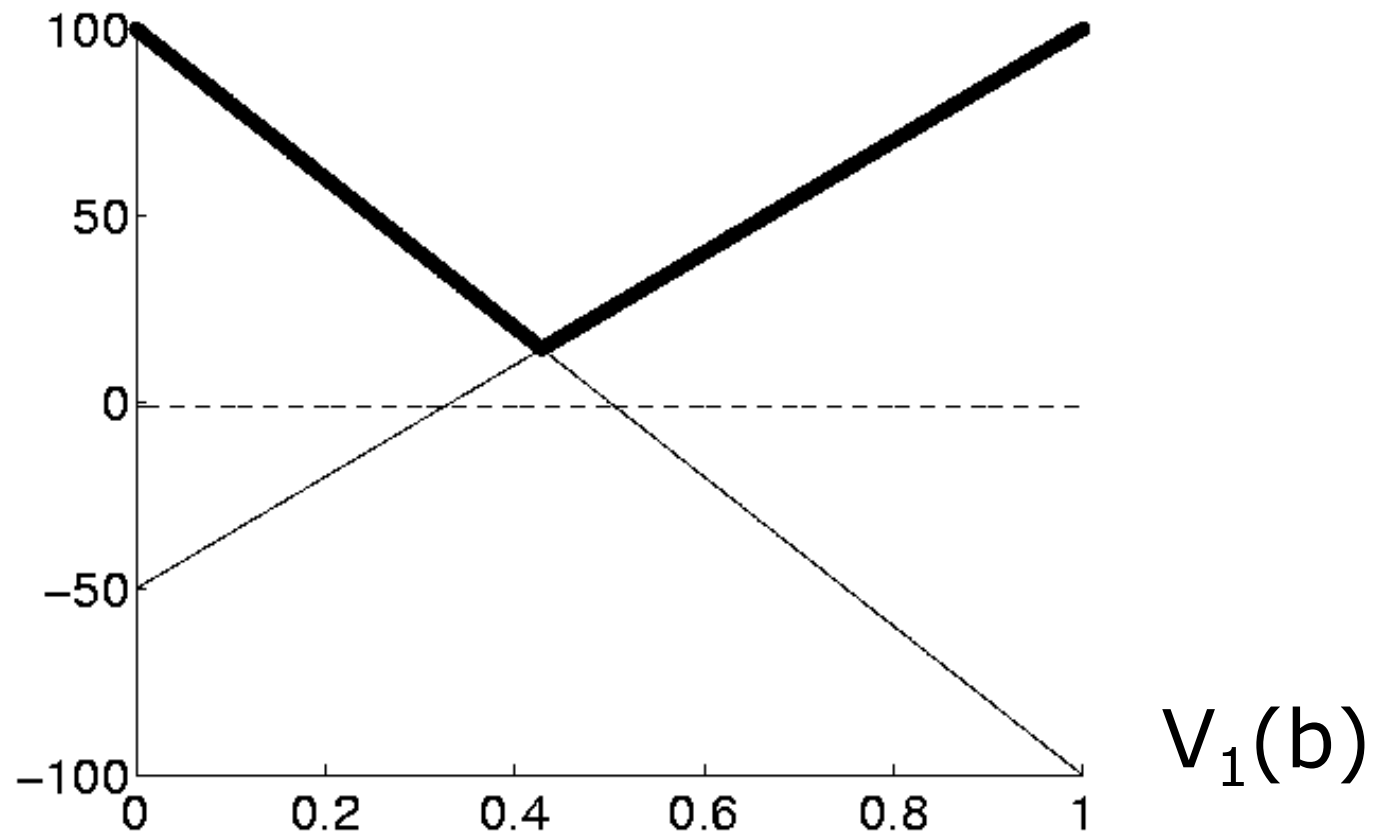
- With  $V_1(b)$ , note that only the first two components contribute.
- The third component can be safely pruned

$$V_1(b) = \max \left\{ \begin{array}{cc} -100 p_1 & +100 (1 - p_1) \\ 100 p_1 & -50 (1 - p_1) \end{array} \right\}$$



# Incorporating Observation

- Suppose that the robot can receive an observation before deciding on an action.



# Incorporating Observation

- Suppose it perceives  $z_1$ :  $p(z_1 | x_1) = 0.7$  and  $p(z_1 | x_2) = 0.3$ .
- Given the obs  $z_1$  we update the belief using **Bayes rule**.

$$p'_1 = \frac{0.7 p_1}{p(z_1)} \quad \text{where} \quad p(z_1) = 0.7 p_1 + 0.3(1 - p_1) = 0.4 p_1 + 0.3$$

- Now,  $V_1(b | z_1)$  is given by

$$\begin{aligned} V_1(b | z_1) &= \max \left\{ \begin{array}{l} -100 \cdot \frac{0.7 p_1}{p(z_1)} + 100 \cdot \frac{0.3 (1-p_1)}{p(z_1)} \\ 100 \cdot \frac{0.7 p_1}{p(z_1)} - 50 \cdot \frac{0.3 (1-p_1)}{p(z_1)} \end{array} \right\} \\ &= \frac{1}{p(z_1)} \max \left\{ \begin{array}{l} -70 p_1 + 30 (1 - p_1) \\ 70 p_1 - 15 (1 - p_1) \end{array} \right\} \end{aligned}$$

# Expected Value after Measuring

- But, we do not know *in advance* what the next measurement will be,
- So we must compute the expected belief

$$\begin{aligned}\bar{V}_1(b) &= E_z[V_1(b | z)] = \sum_{i=1}^2 p(z_i) V_1(b | z_i) \\ &= \sum_{i=1}^2 p(z_i) V_1\left(\frac{p(z_i | x_1) p_1}{p(z_i)}\right) \\ &= \sum_{i=1}^2 V_1(p(z_i | x_1) p_1)\end{aligned}$$

# Expected Value after Measuring

- But, we do not know *in advance* what the next measurement will be,
- So we must compute the expected belief

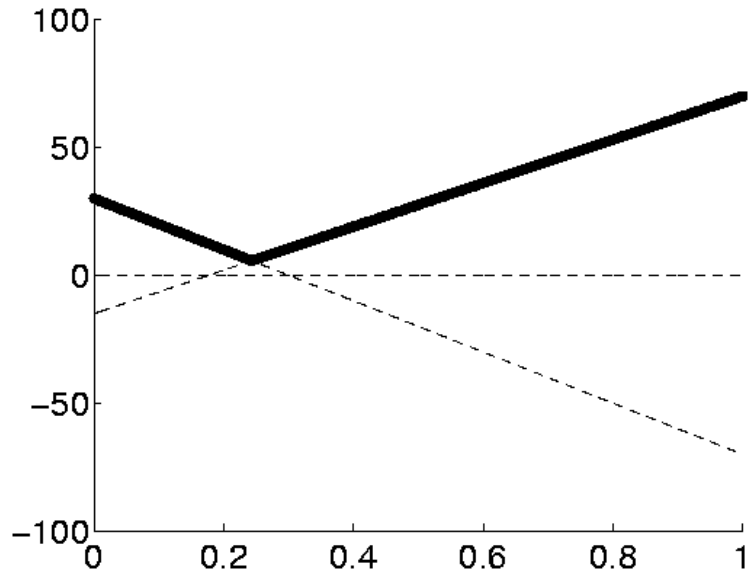
$$\begin{aligned}\bar{V}_1(b) &= E_z[V_1(b | z)] \\ &= \sum_{i=1}^2 p(z_i) V_1(b | z_i) \\ &= \max \left\{ \begin{array}{cc} -70 p_1 & +30 (1 - p_1) \\ 70 p_1 & -15 (1 - p_1) \end{array} \right\} \\ &\quad + \max \left\{ \begin{array}{cc} -30 p_1 & +70 (1 - p_1) \\ 30 p_1 & -35 (1 - p_1) \end{array} \right\}\end{aligned}$$

# Resulting Value Function

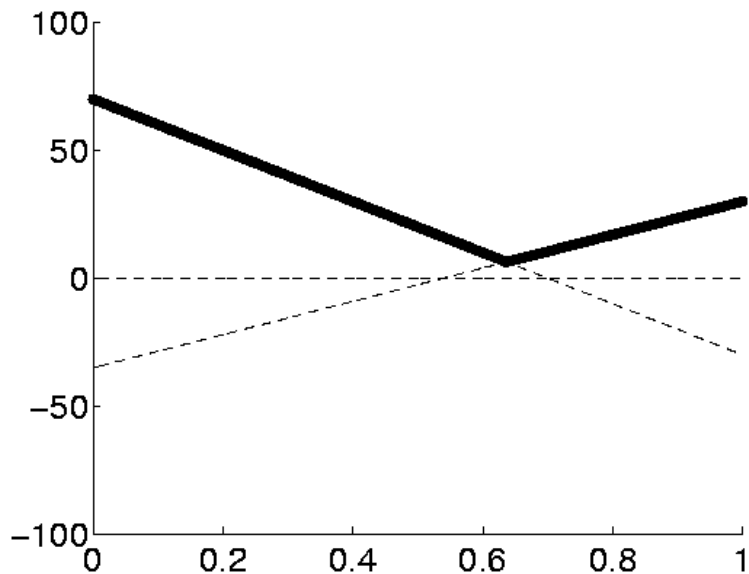
- The four possible combinations yield the following function which then can be simplified and pruned.

$$\begin{aligned}\bar{V}_1(b) &= \max \left\{ \begin{array}{cccc} -70 p_1 & +30 (1 - p_1) & -30 p_1 & +70 (1 - p_1) \\ -70 p_1 & +30 (1 - p_1) & +30 p_1 & -35 (1 - p_1) \\ +70 p_1 & -15 (1 - p_1) & -30 p_1 & +70 (1 - p_1) \\ +70 p_1 & -15 (1 - p_1) & +30 p_1 & -35 (1 - p_1) \end{array} \right\} \\ &= \max \left\{ \begin{array}{cc} -100 p_1 & +100 (1 - p_1) \\ +40 p_1 & +55 (1 - p_1) \\ +100 p_1 & -50 (1 - p_1) \end{array} \right\}\end{aligned}$$

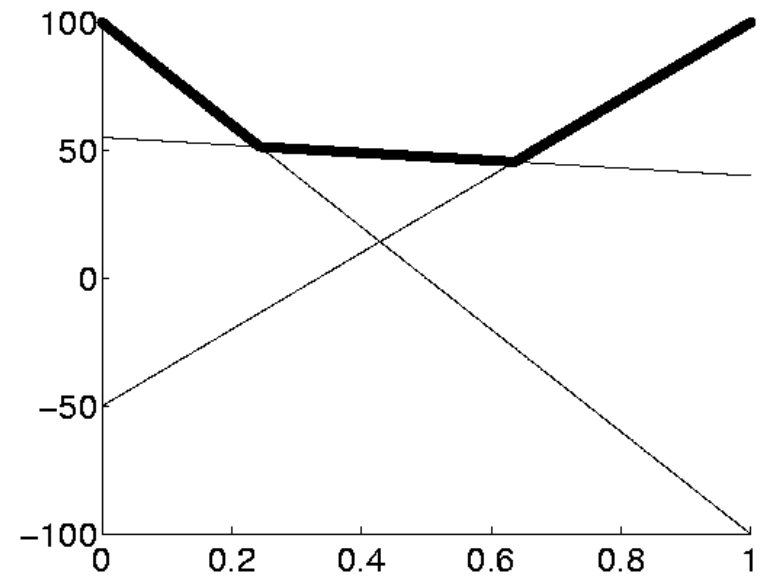
# Value Function



$$p(z_1) V_1(b|z_1)$$



$$p(z_2) V_2(b|z_2)$$



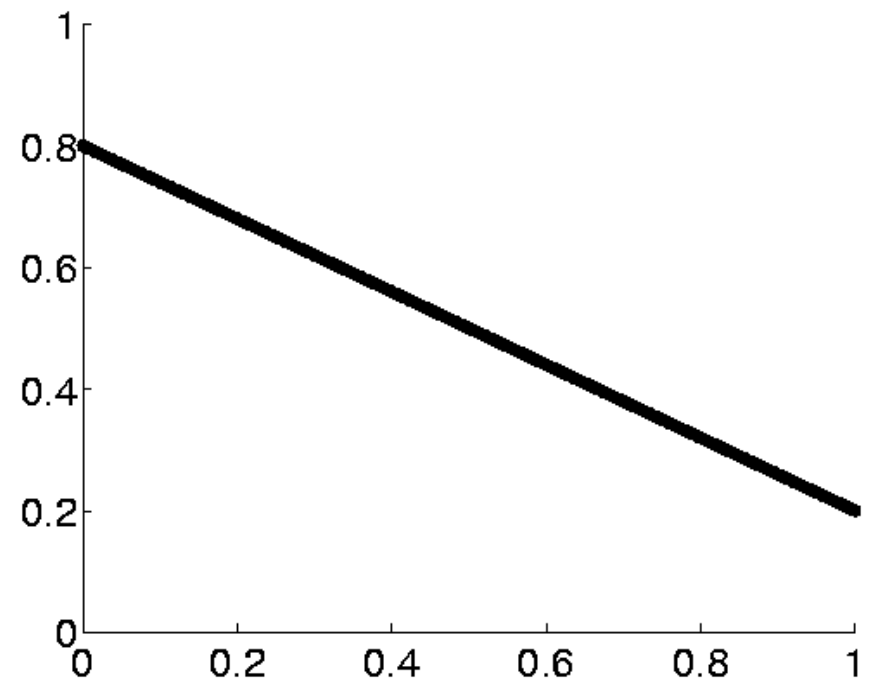
$$\bar{V}_1(b)$$

# Increasing the Time Horizon

- When the agent selects  $u_3$  its state may change.
- When computing the value function, we have to take these potential state changes into account.

$$\begin{aligned} p'_1 &= E_x[p(x_1 | x, u_3)] \\ &= \sum_{i=1}^2 p(x_1 | x_i, u_3)p_i \\ &= 0.2p_1 + 0.8(1 - p_1) \\ &= 0.8 - 0.6p_1 \end{aligned}$$

$P(x=x_1 \text{ after executing } u_3)$



$P(x=x_1 \text{ originally})$

# Resulting Value Function after executing $u_3$

Taking the state transitions into account, we finally obtain.

$$\bar{V}_1(b) = \max \left\{ \begin{array}{cccc} -70 p_1 & +30 (1 - p_1) & -30 p_1 & +70 (1 - p_1) \\ -70 p_1 & +30 (1 - p_1) & +30 p_1 & -35 (1 - p_1) \\ +70 p_1 & -15 (1 - p_1) & -30 p_1 & +70 (1 - p_1) \\ +70 p_1 & -15 (1 - p_1) & +30 p_1 & -35 (1 - p_1) \end{array} \right\}$$

$$= \max \left\{ \begin{array}{cc} -100 p_1 & +100 (1 - p_1) \\ +40 p_1 & +55 (1 - p_1) \\ +100 p_1 & -50 (1 - p_1) \end{array} \right\}$$

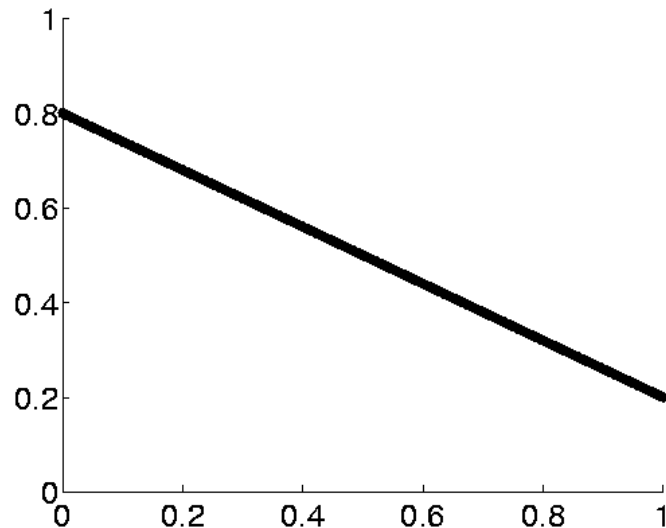
$$\bar{V}_1(b | u_3) = \max \left\{ \begin{array}{cc} 60 p_1 & -60 (1 - p_1) \\ 52 p_1 & +43 (1 - p_1) \\ -20 p_1 & +70 (1 - p_1) \end{array} \right\}$$



# Value Function after executing $u_3$

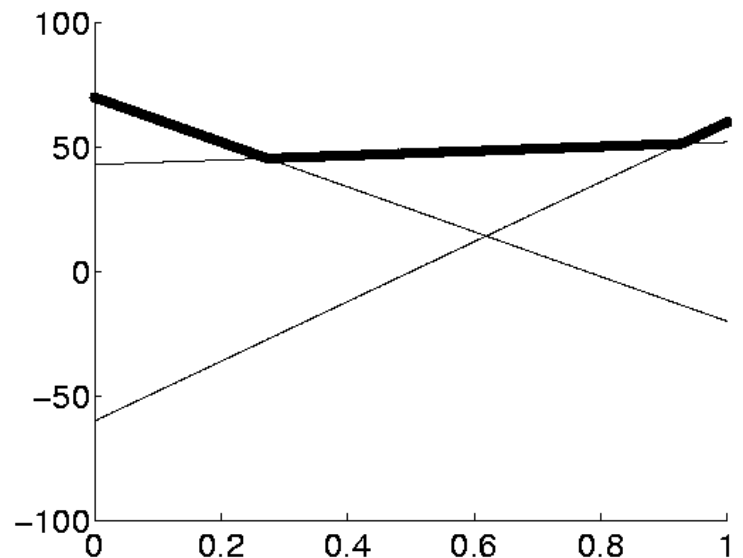
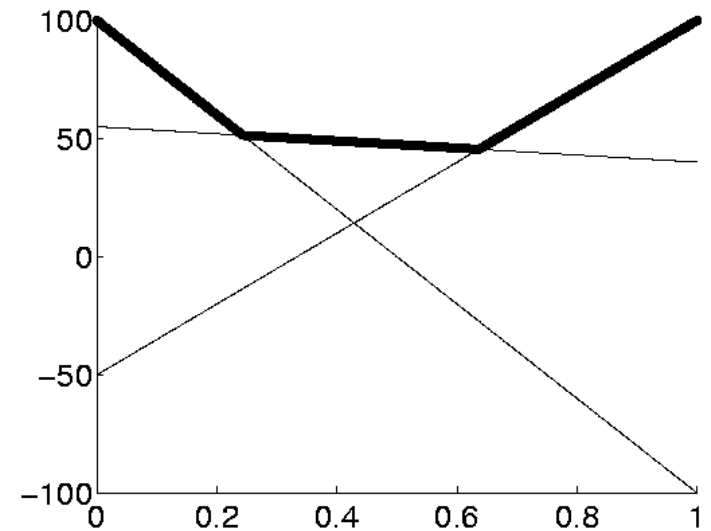
$$\bar{V}_1(b)$$

$P(x=x_1 \text{ after executing } u_3)$



$P(x=x_1 \text{ originally})$

$$\bar{V}_1(b|u_3)$$

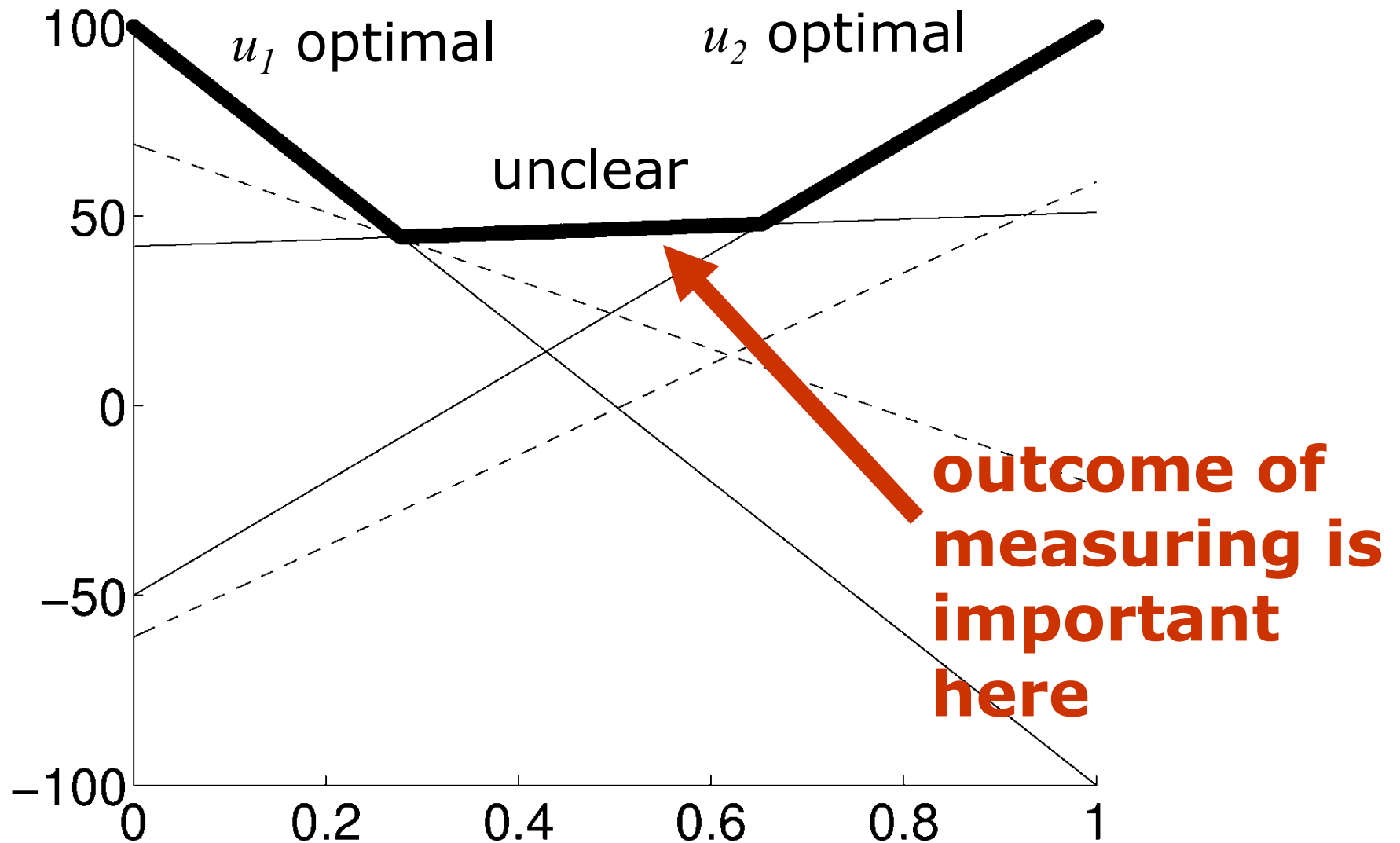


# Value Function for T=2

- Taking into account that the agent can either directly perform  $u_1$  or  $u_2$  or first  $u_3$  and then  $u_1$  or  $u_2$ , we obtain (after pruning)

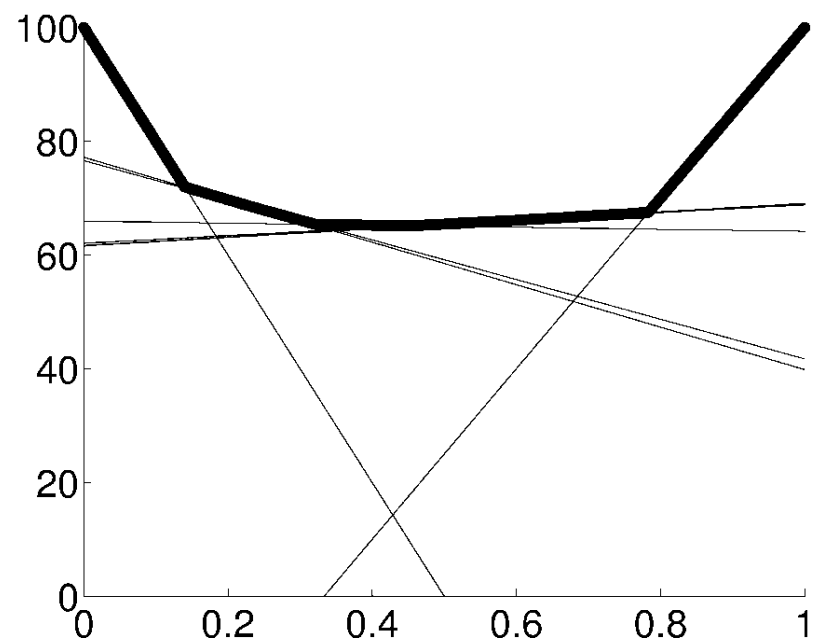
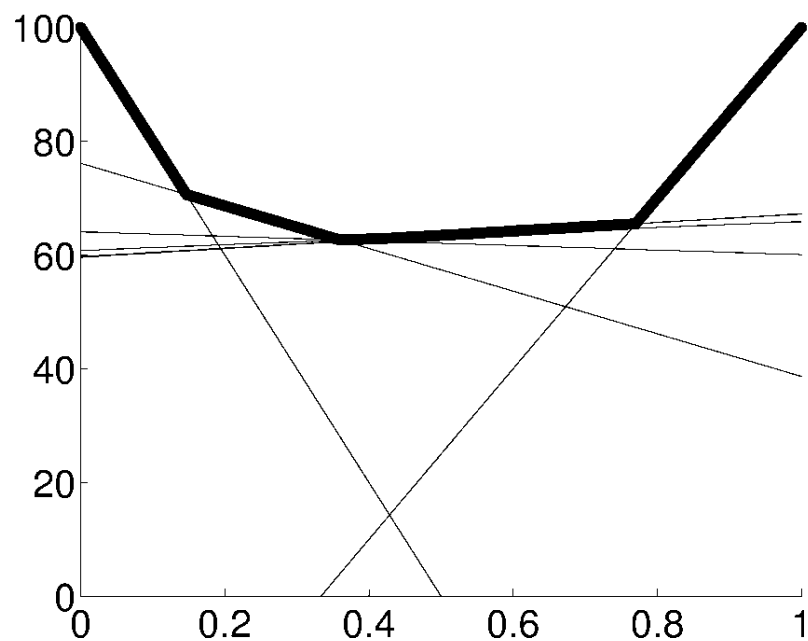
$$\bar{V}_2(b) = \max \left\{ \begin{array}{ll} -100 p_1 & +100 (1 - p_1) \\ 100 p_1 & -50 (1 - p_1) \\ 51 p_1 & +42 (1 - p_1) \end{array} \right\}$$

# Graphical Representation of $V_2(b)$

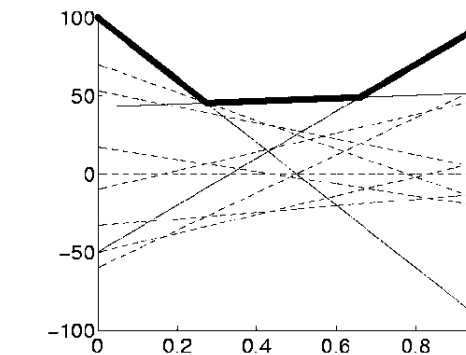
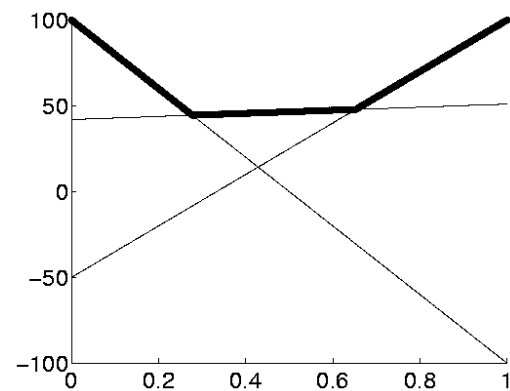
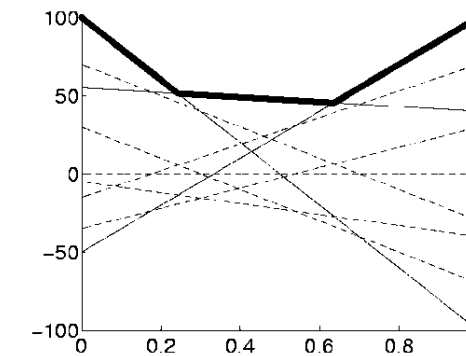
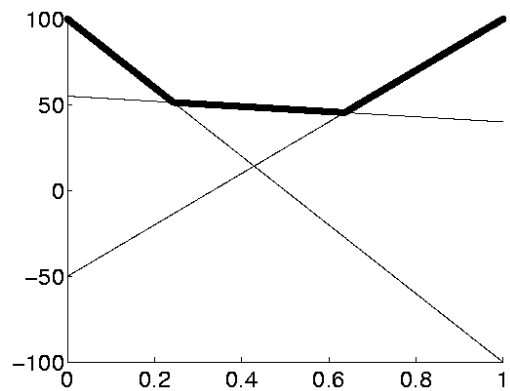
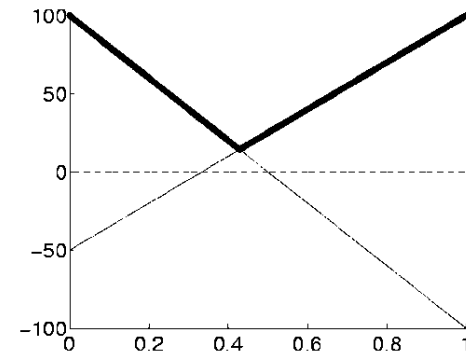
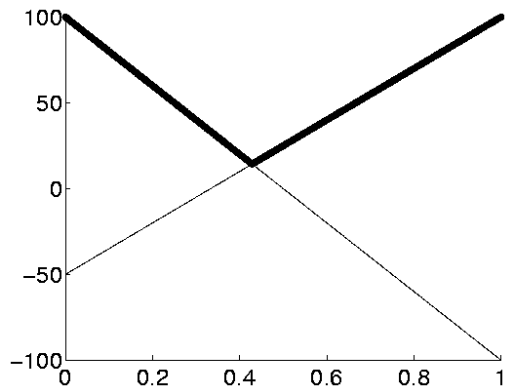


# Deep Horizons

- We have now completed a full backup in belief space.
- This process can be applied recursively.
- The value functions for  $T=10$  and  $T=20$  are



# Deep Horizons and Pruning



# Why Pruning is Essential

- Each **update introduces additional linear components** to  $V$ .
- Each **measurement squares the number of linear components**.
- Thus, an unpruned value function for  $T=20$  includes more than  $10^{547,864}$  linear functions.
- At  $T=30$  we have  $10^{561,012,337}$  linear functions.
- The pruned value functions at  $T=20$ , in comparison, contains only 12 linear components.
- The combinatorial explosion of linear components in the value function are the major reason why **exact solution of POMDPs is usually impractical**

# POMDP Approximations

- Point-based value iteration
- QMDPs
- AMDPs

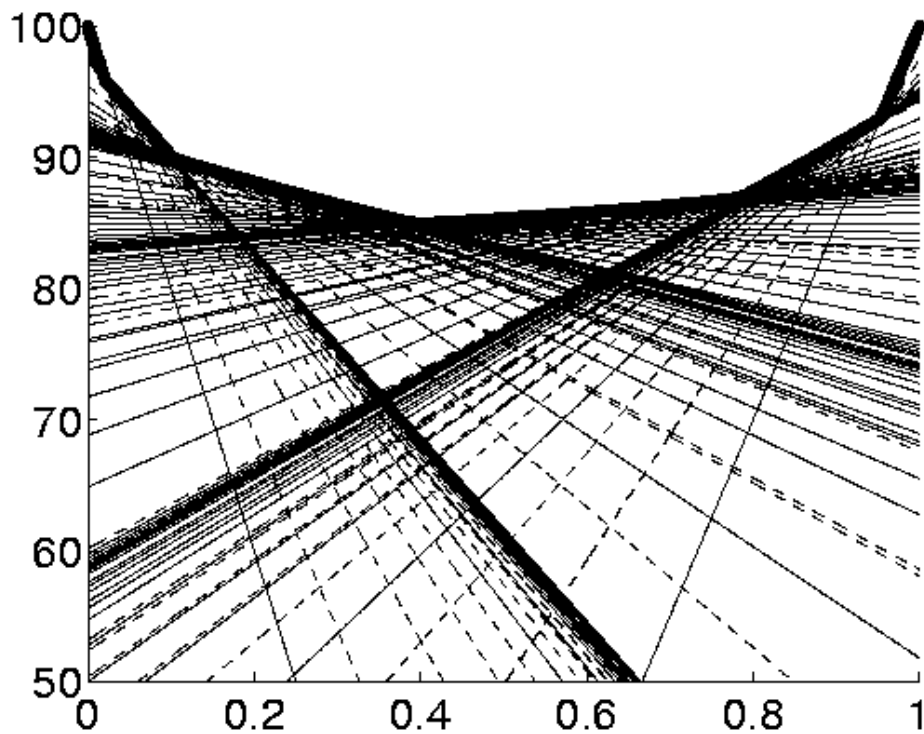
# Point-based Value Iteration

- Maintains a set of example beliefs
- Only considers constraints that maximize value function for at least one of the examples

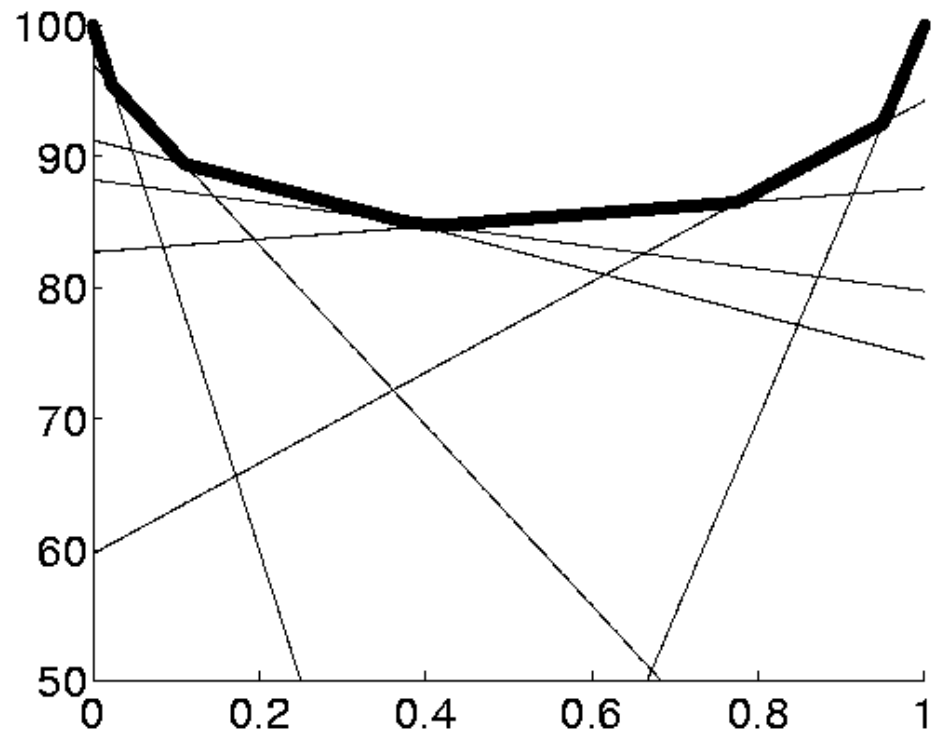


# Point-based Value Iteration

Value functions for  $T=30$



Exact value function



PBVI

# QMDPs

- QMDPs only consider state uncertainty in the first step
- After that, assume that the world is fully observable.

# POMDP Summary

- POMDPs compute the optimal action in partially observable, stochastic domains.
- For finite horizon problems, the resulting value functions are piecewise linear and convex.
- In each iteration the number of linear constraints grows exponentially.
- Until recently, POMDPs only applied to very small state spaces with small numbers of possible observations and actions.
  - But with PBVI,  $|S| = \text{millions}$