# Parameter-Independent Strategies for pMDPs via POMDPs

Sebastian Arming — UNIVERSITÄT SALZBURG

Ezio Bartocci — TU WIEN TECHNISCHE UNIVERSITÄT WIEN

Krishnendu Chatterjee — IST AUSTRIA Institute of Science and Technology

Joost-Pieter Katoen — RWTH AACHEN UNIVERSITY

Ana Sokolova — UNIVERSITÄT SALZBURG

RiSE
Rigorous Systems Engineering

# The problem

Finding **policies**

of a **parametric MDP**

that are **optimal**

over the **whole parameter space**.

# The problem

Finding **policies**

of a **parametric MDP**

that are **expectation optimal**

(over the **whole parameter space**).

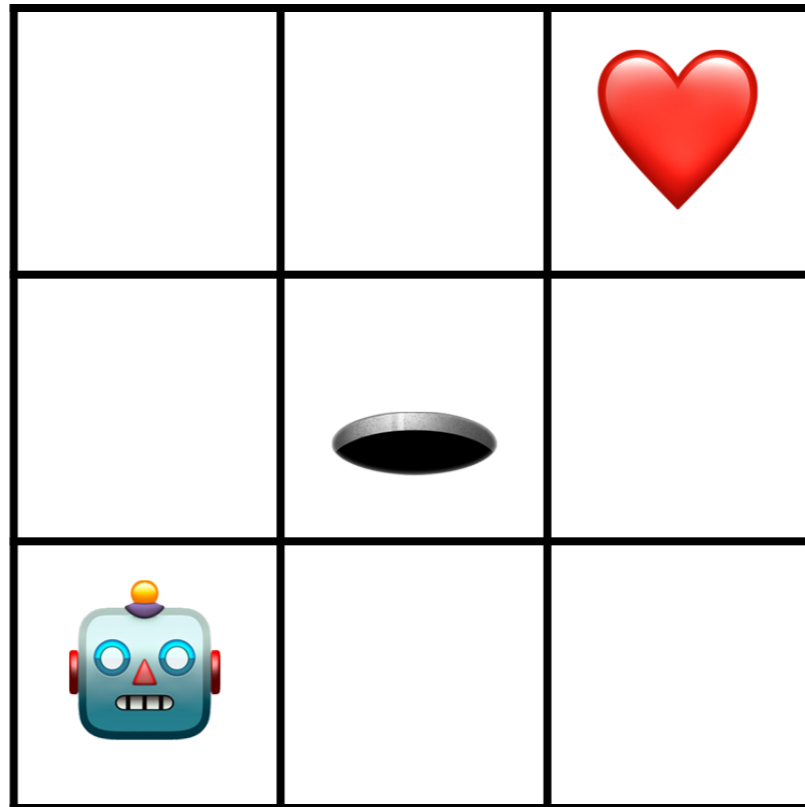# The solution

Finding **policies**

of a **parametric MDP**

that are **expectation optimal**
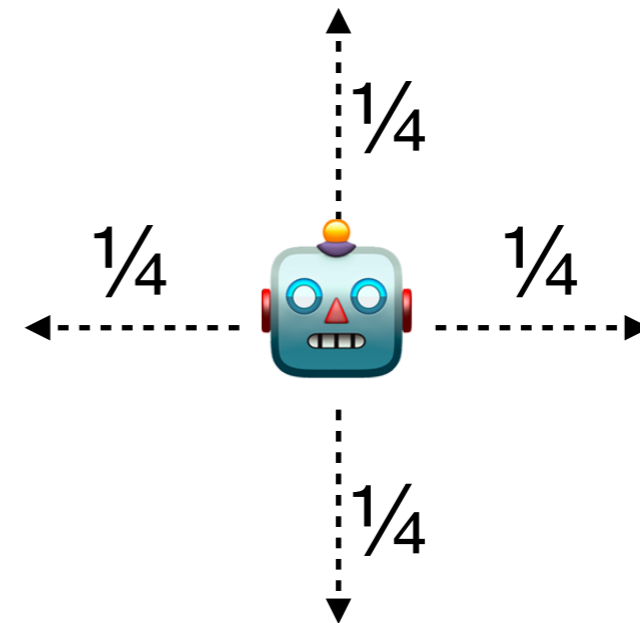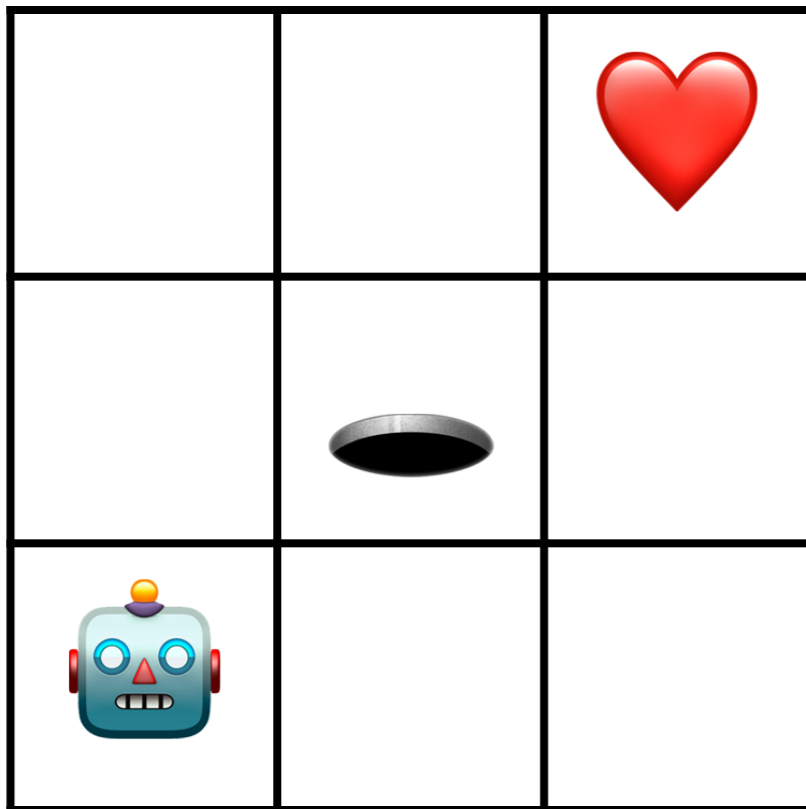
(over the **whole parameter space**)

amounts to solving a **suitable POMDP.**

# MDP: Markov Decision Process

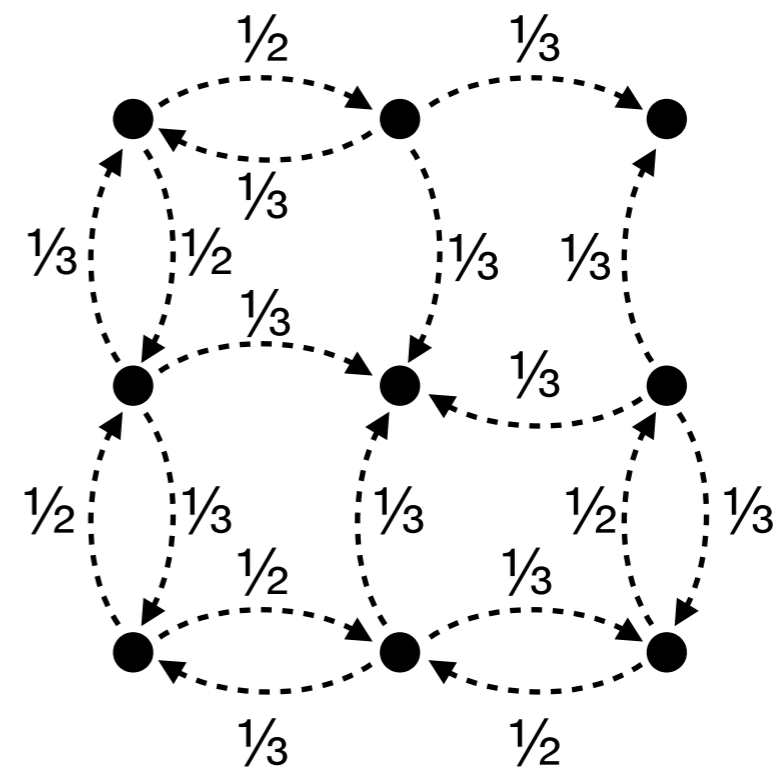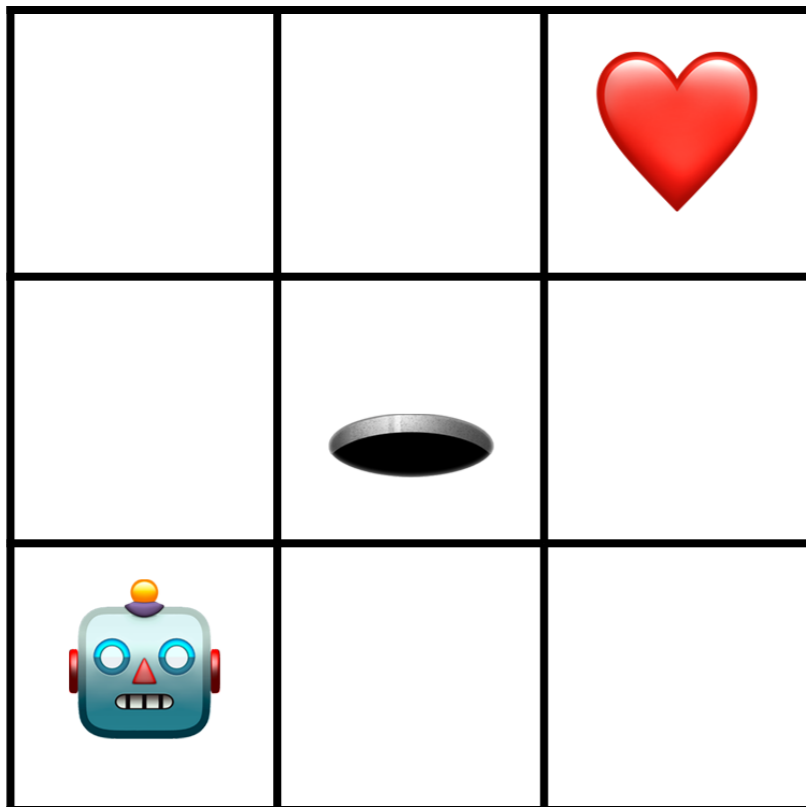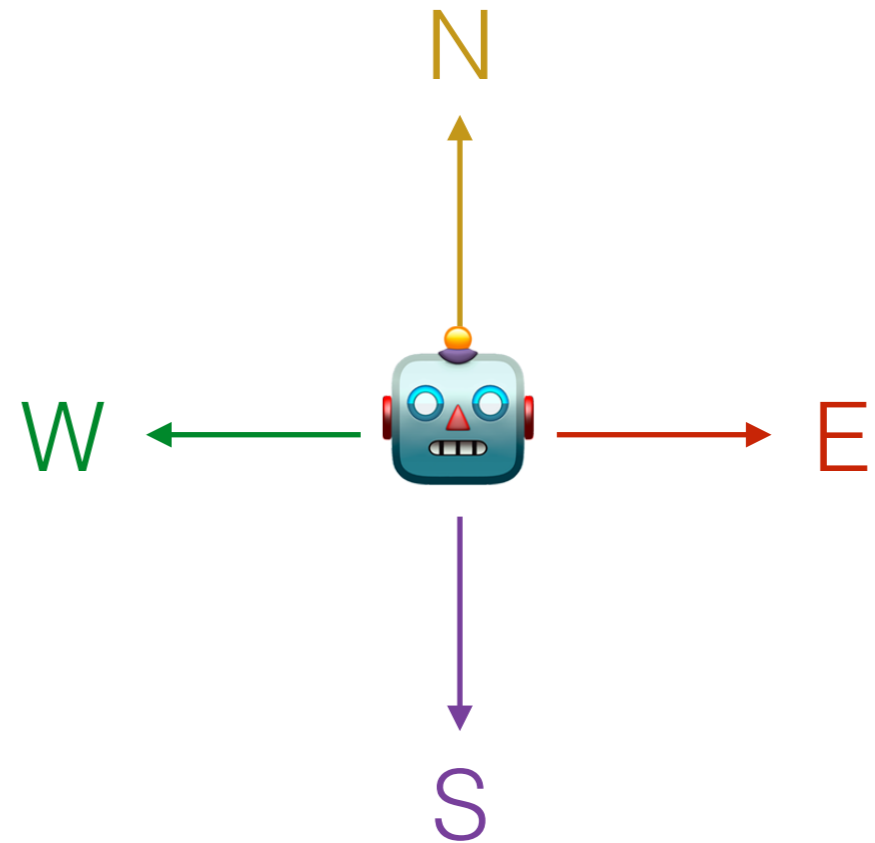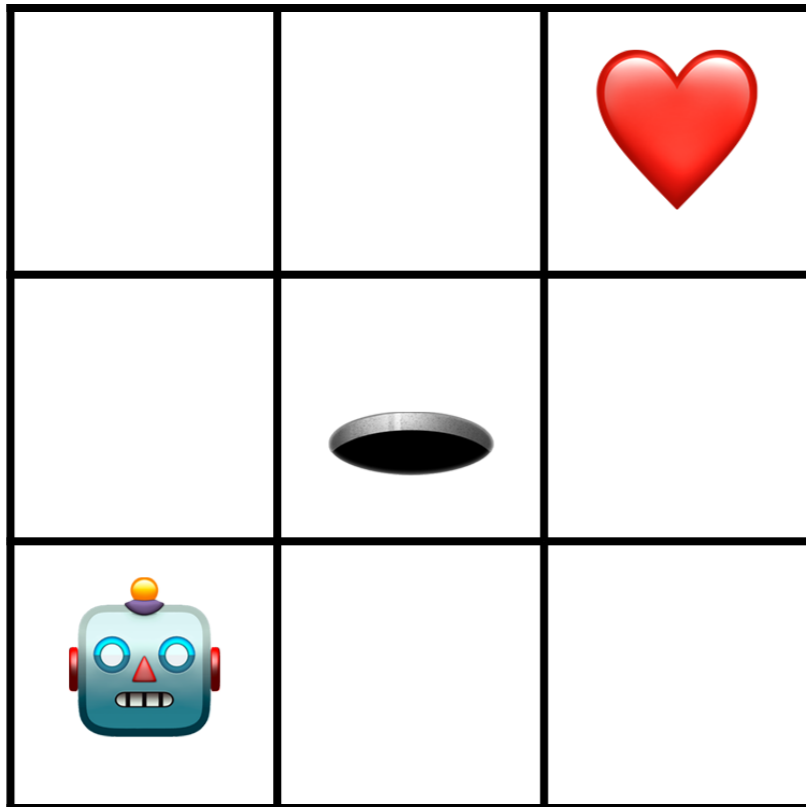# Robot example

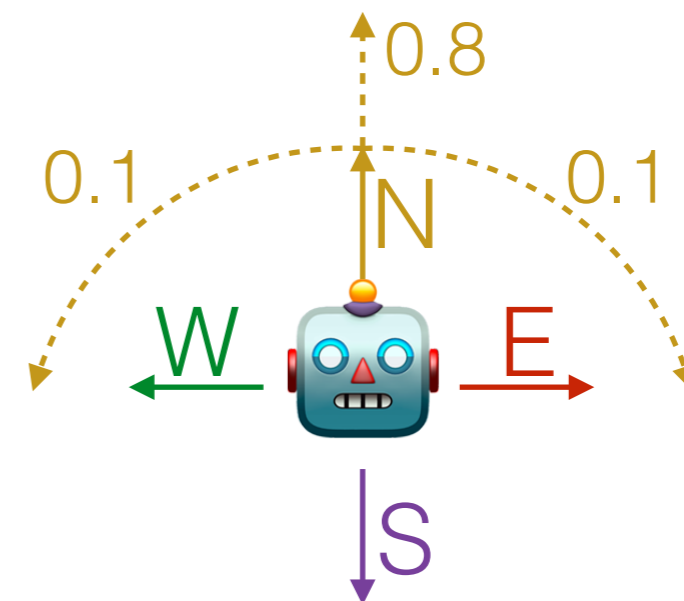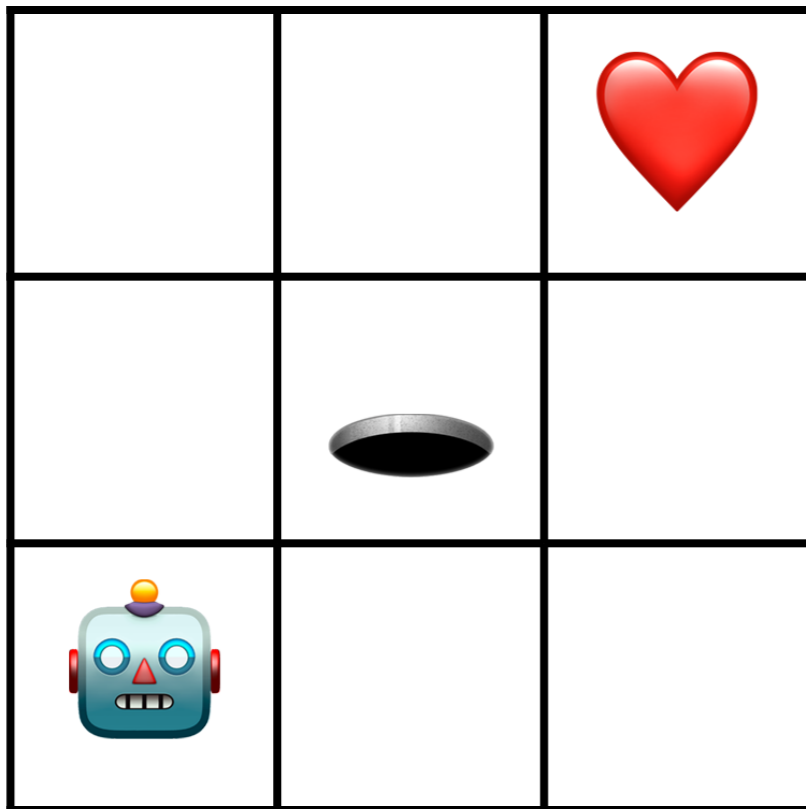# Robot example

# Markov Chain

$$\mathbf{Pr}(\heartsuit) = \frac{1}{7}$$

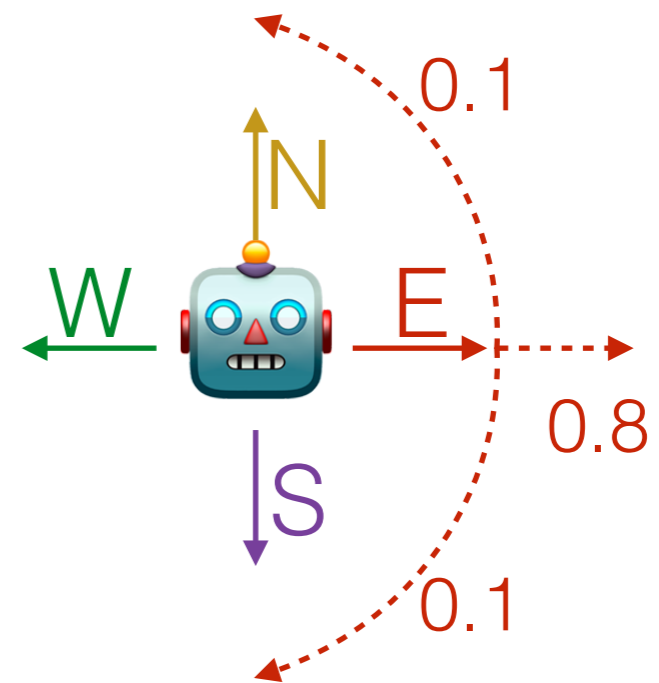A **MC** is a pair *(S,T)* where
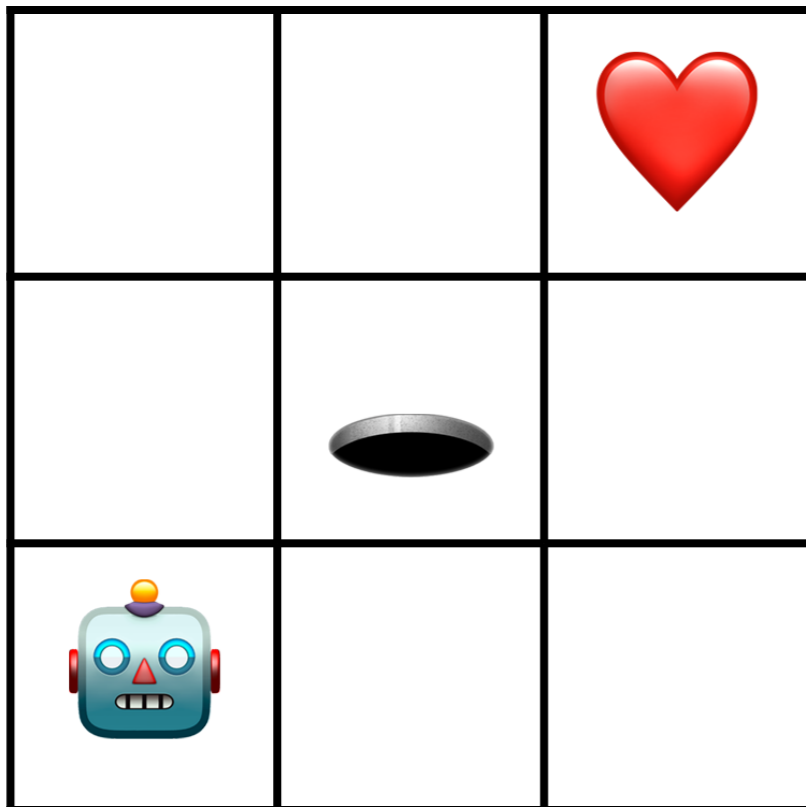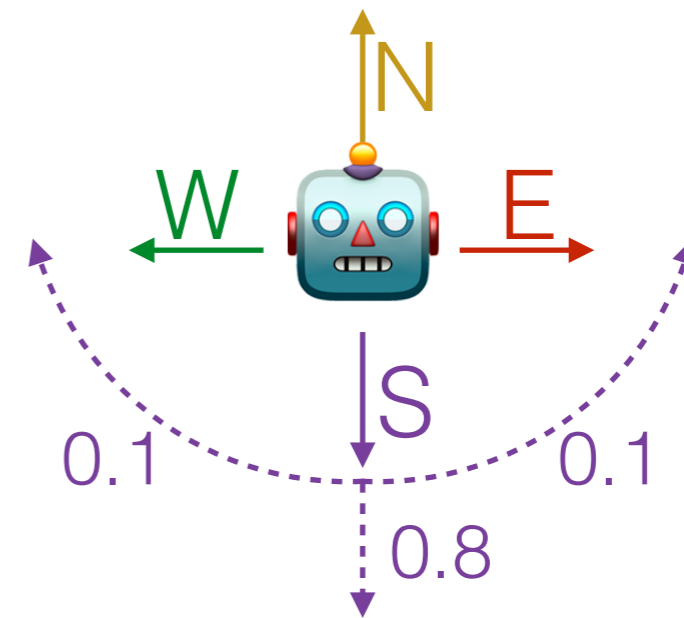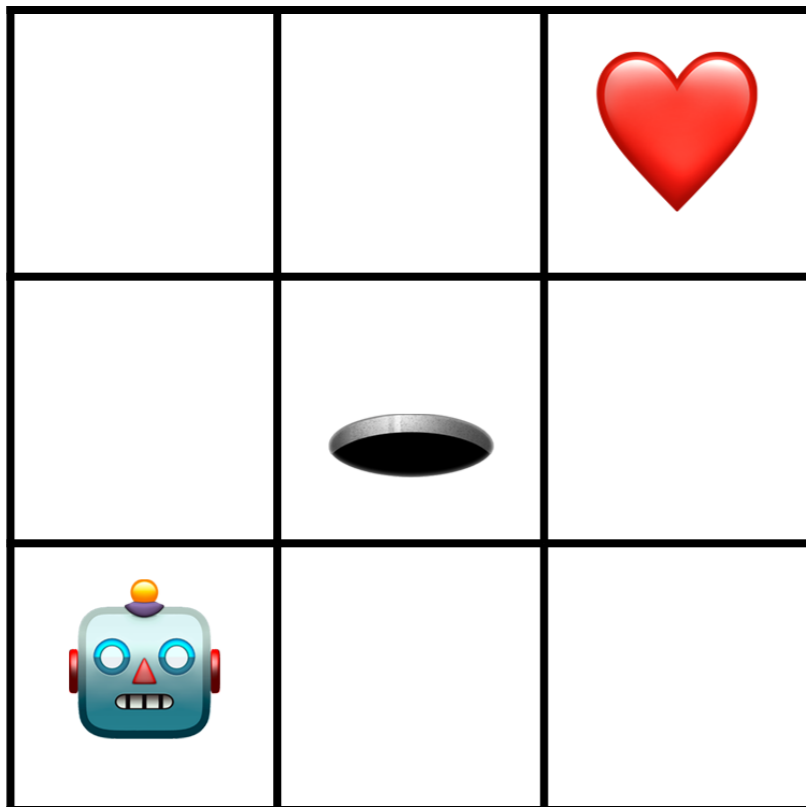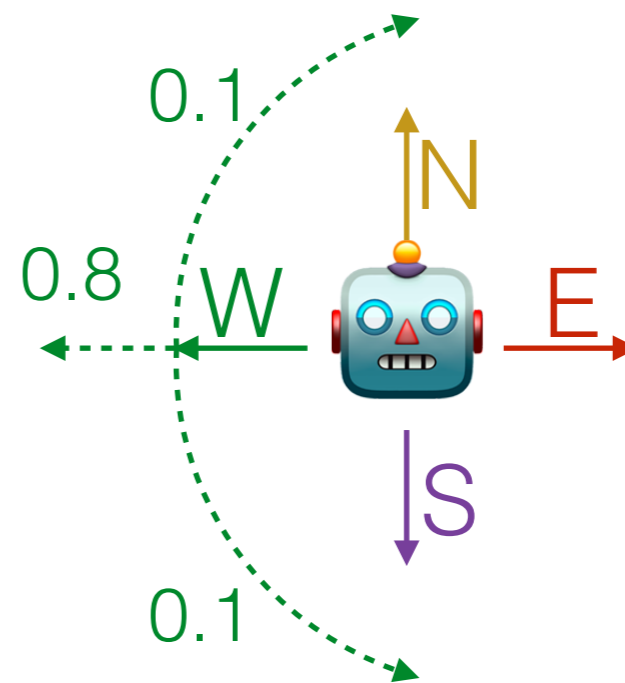
- *S* is a set of **states**
- *T: S → 𝒟S* is a **transition function**
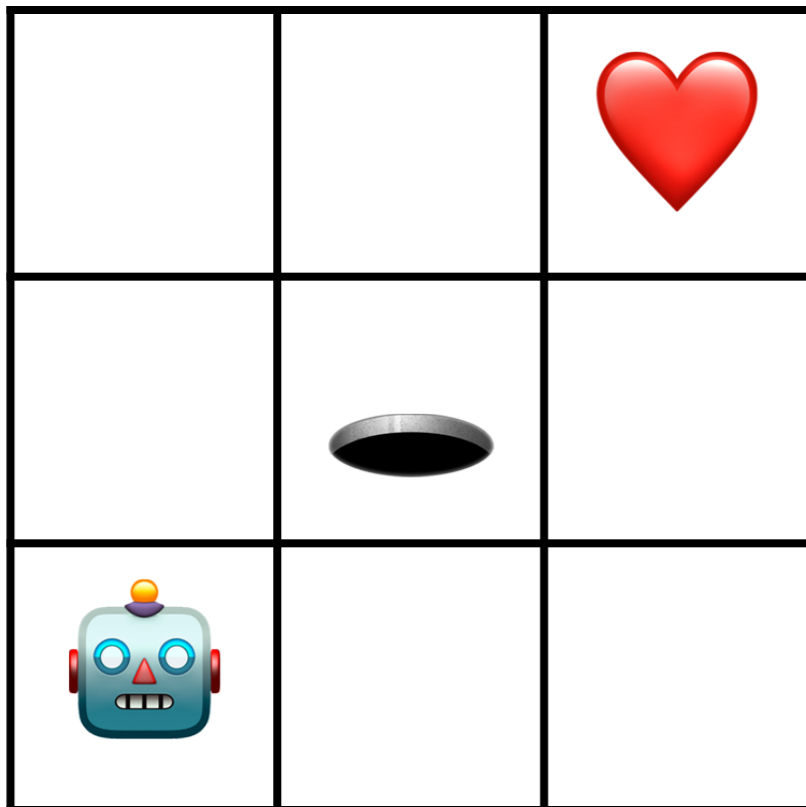
# Robot example

# Robot example

# Robot example

# Robot example

# Robot example

# Markov Decision Process

# Markov Decision Process MDP

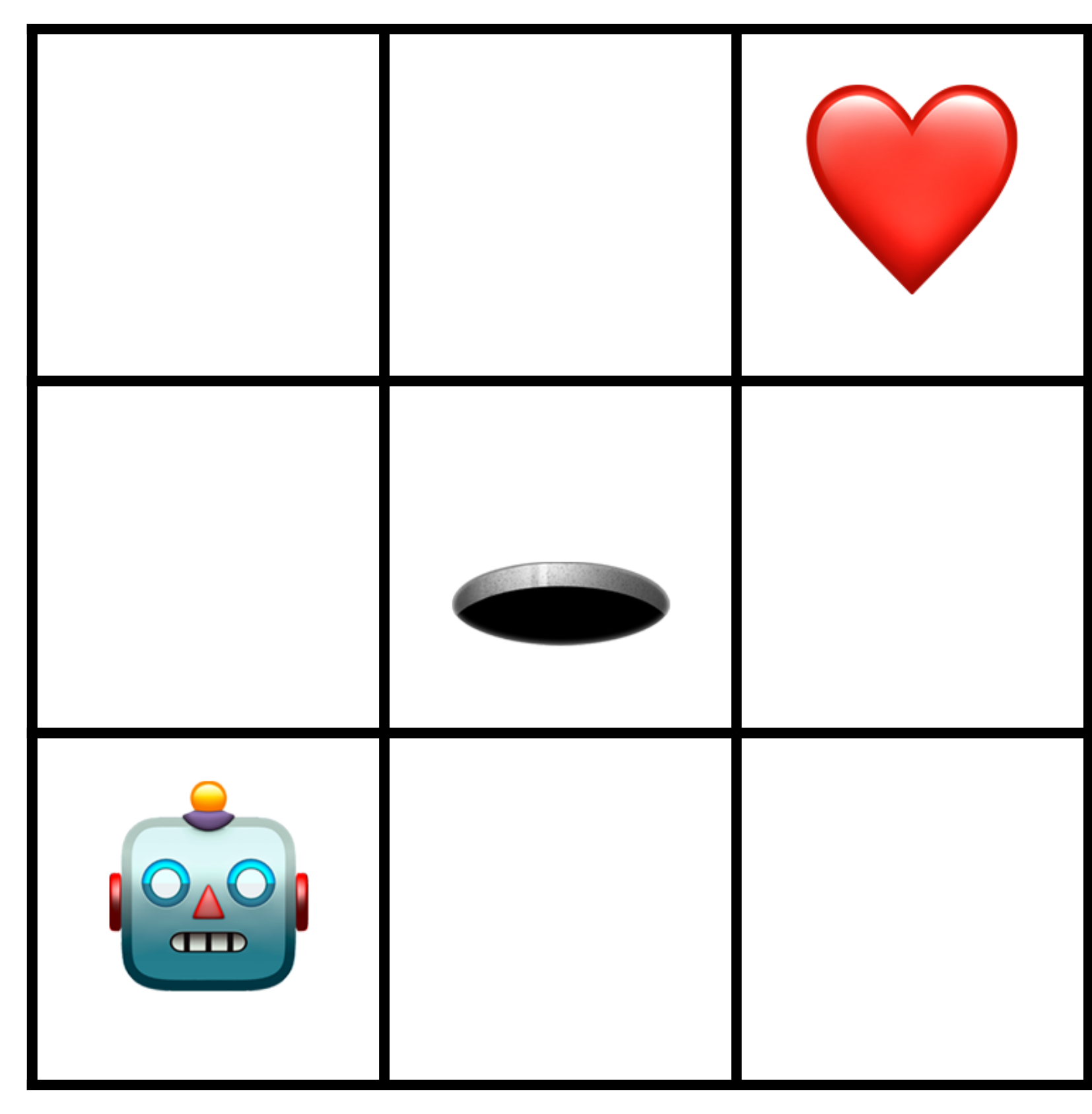

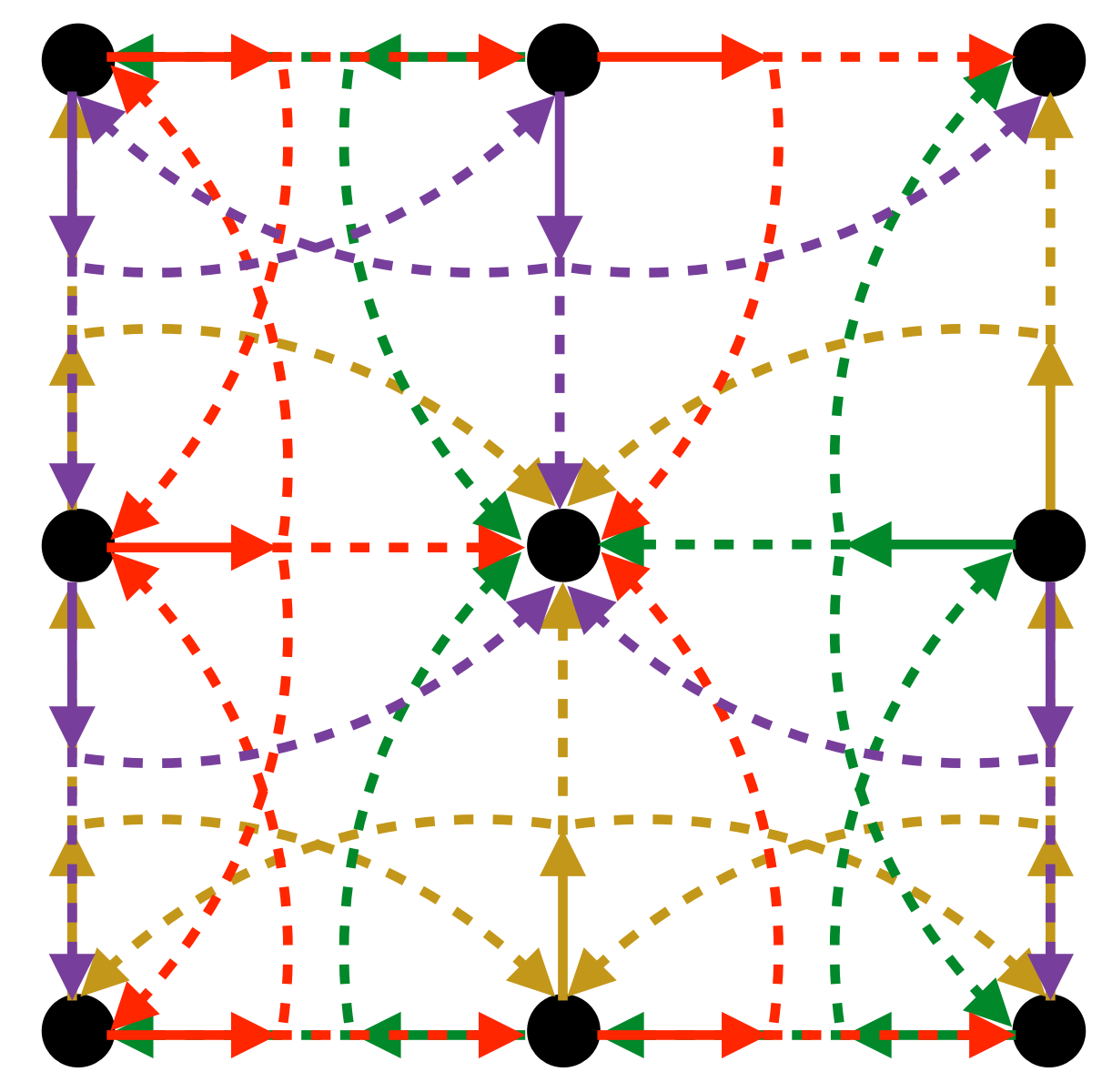

A **MDP** is a tuple $(S, A, T)$ where

- $S$ is a set of **states**
- $A$ is a set of **actions**
- $T: S \times A \to \mathscr{D}S$ is a **transition function**

# Markov Decision Process

$$\mathbf{Pr}(\text{❤}) = f(\text{🤖})$$

# Markov Decision Process



$$\textbf{Pr}(\heartsuit) = f(\text{🤖})$$

# Markov Decision Process



**Pr**(❤️)= f(🤖)

# Policy

A **policy** $\pi$ for an MDP *(S,A,T)* is a function

$$\pi: (S \times A)_* \times S \rightarrow \mathscr{D}A$$

**deterministic:** only Dirac distributions     $\pi: (S \times A)_* \times S \rightarrow A$
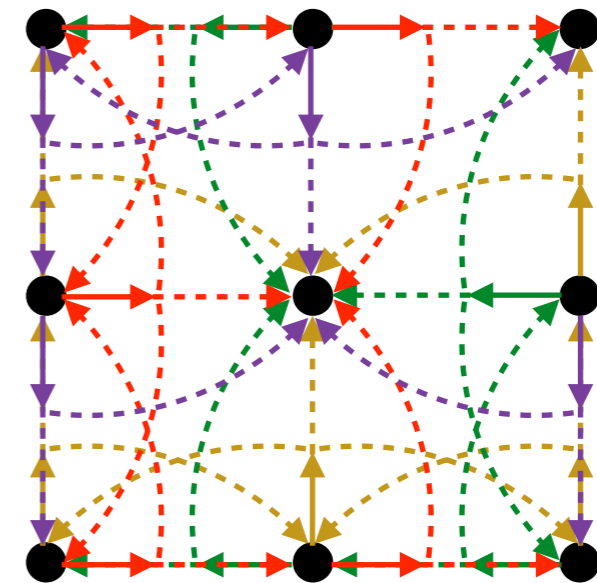
**memoryless:** $\pi(\ldots s)=\pi(s)$     $\pi: S \rightarrow \mathscr{D}A$

**simple:** deterministic & memoryless     $\pi: S \rightarrow A$

**Wanted:** a policy that optimizes an objective, e.g. reachability

# Objective

A policy $\pi$ and initial distribution $i$ gives us a probability space

$$(\text{Runs, Cones, } \mathbf{P}_{\pi,i})$$

in the usual way.

Runs — all infinite runs in $(S \times A)^{\omega}$

Cones — the $\sigma$ algebra generated by the
      sets of runs with a common finite prefix (history)

$\mathbf{P}_{\pi,i}$     — the usual measure on cones

A **Borel objective** is a measurable function

$$r: \text{Runs} \rightarrow \mathbb{R}$$

# Objective

A **Borel objective** is a measurable function

$$r: \text{Runs} \to \mathbb{R}$$

Rewards R: $S \times A \to \mathbb{R}$ induce Borel objectives via

$$r_R(s_0, a_0, s_1, a_1, \ldots) = \sum_{i \geq 0} R(s_i, a_i)$$

Reachability objectives are a special case:

Reachability probability = Expectation of reachability objective

# Optimal policy

An **optimal policy** is a policy $\pi$ with

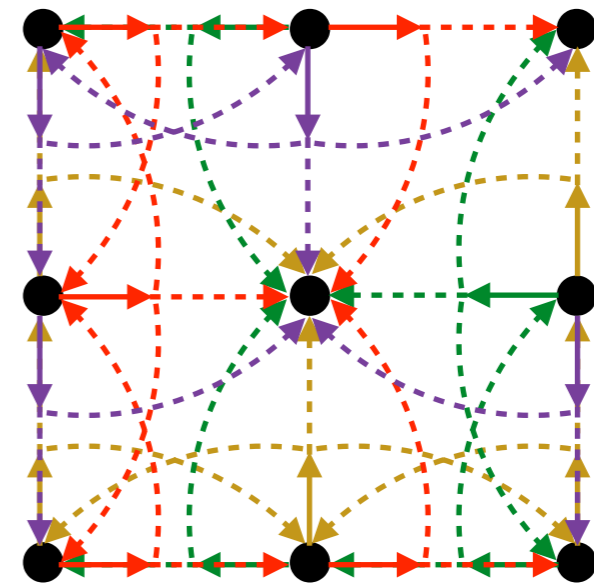$$\mathbf{E}_{\pi,i}(r) = \sup_\sigma \mathbf{E}_{\sigma,i}(r)$$
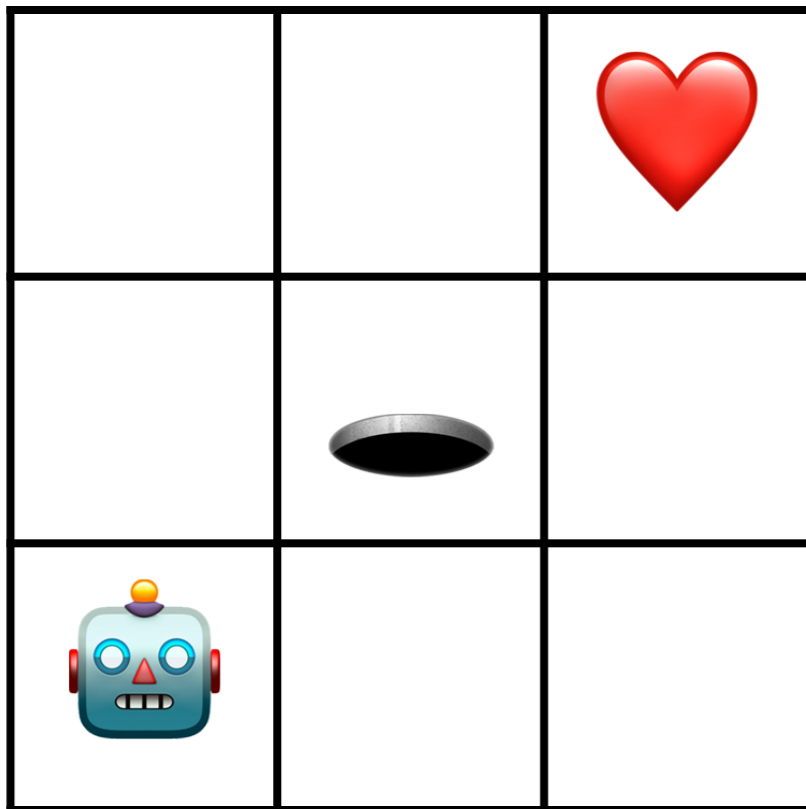
The **expectation** of $\pi$

The **value** of r

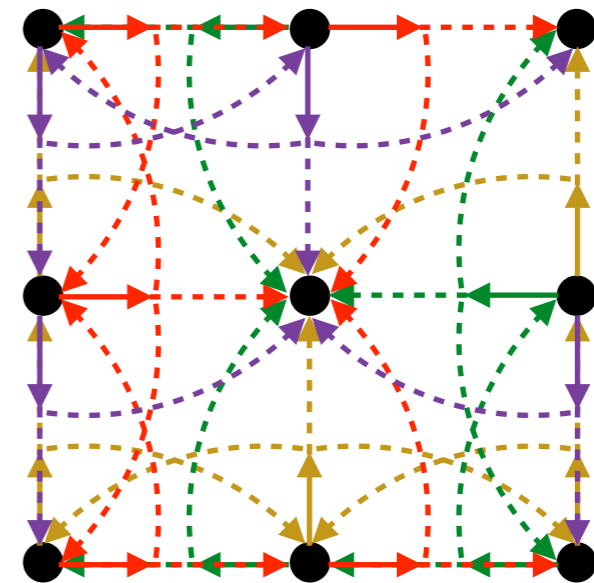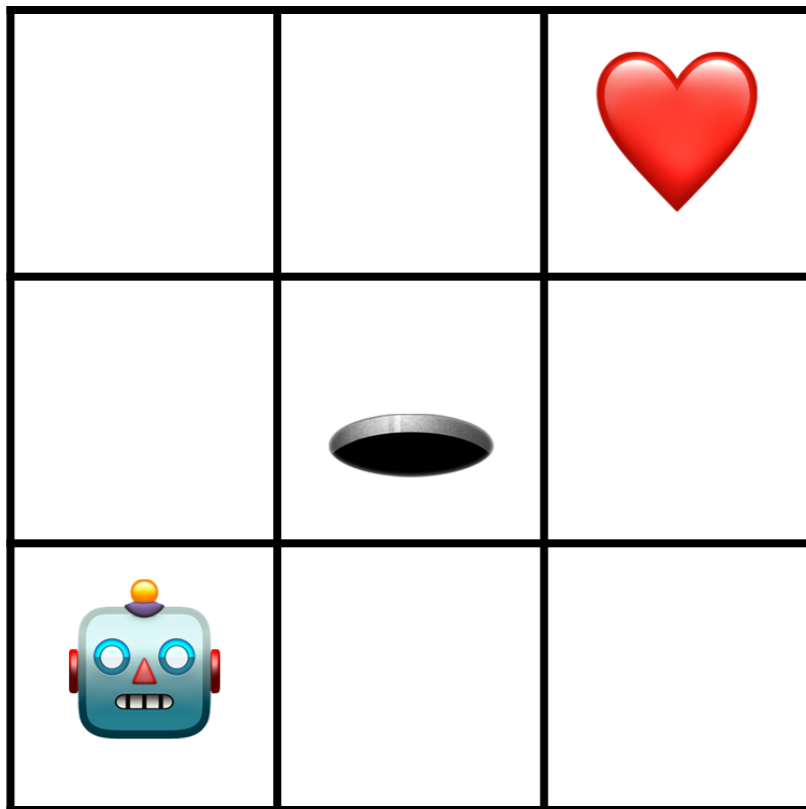An optimal policy not always exists, but **ε-optimal** do:

A policy is **ε-optimal if** $\mathbf{E}_{\pi,i}(r)$ is ε-close to $\sup_\sigma \mathbf{E}_{\sigma,i}(r)$

# Markov Decision Process MDP



$$\mathbf{Pr}(\text{❤}) = 0, \frac{3}{19}, \frac{4}{15}, \frac{3}{13}, \frac{3}{25}, \frac{9}{40}, \frac{1}{8}, \frac{18}{55}, \frac{9}{52}, \frac{9}{100}, \frac{27}{95}$$
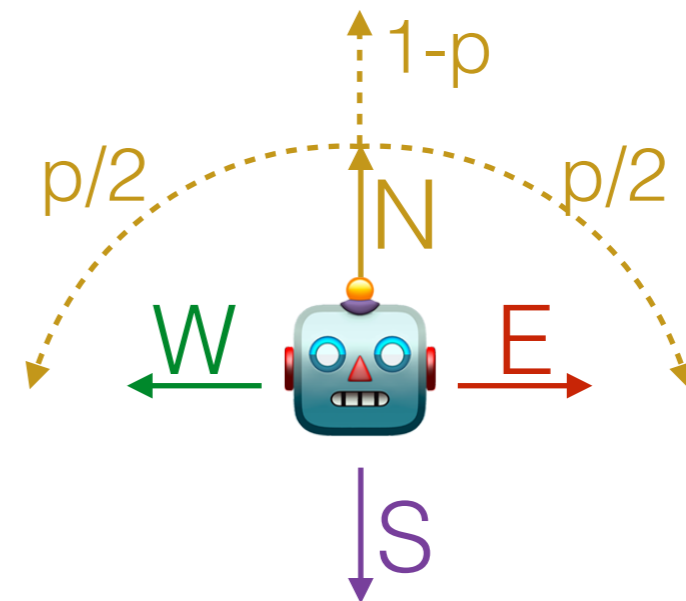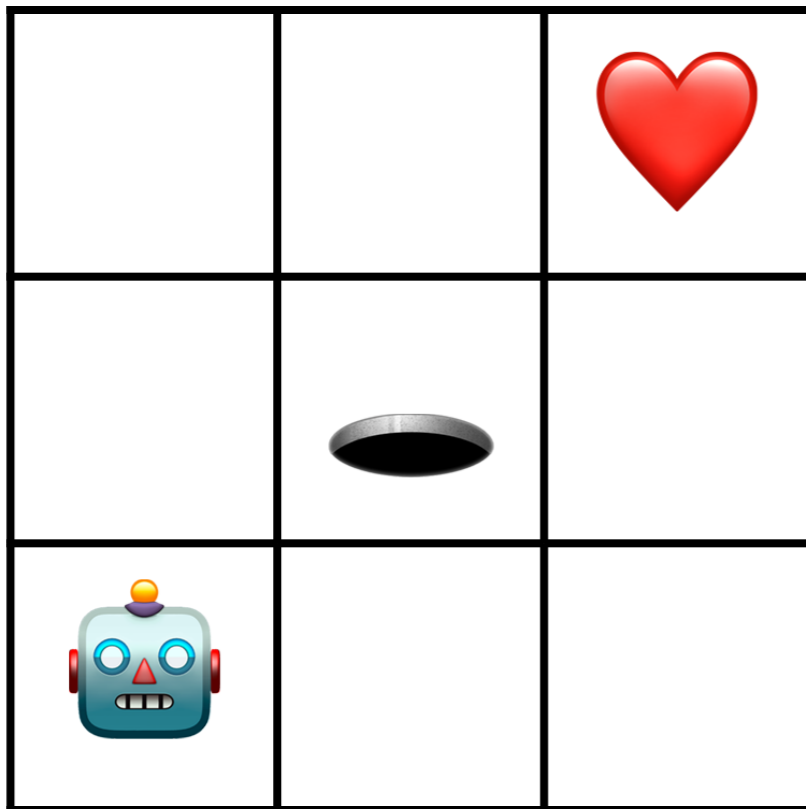
23

# Markov Decision Process <sub>MDP</sub>



$$\mathbf{Pr_{sup}}(\heartsuit) = 18/55$$

# Parameters

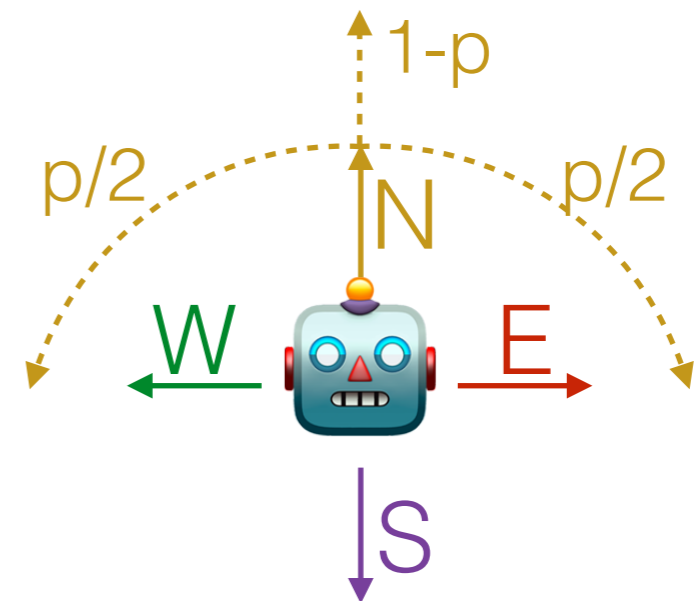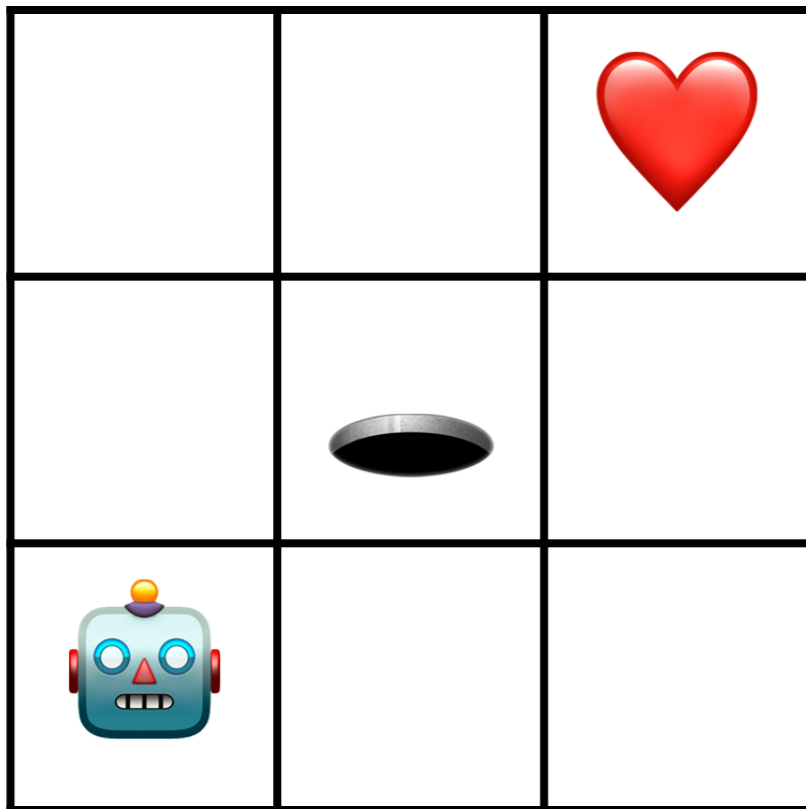# pMDP: parametric MDP

# Parametric Markov Decision Process <span style="font-size:small">pMDP</span>



A **pMDP** is a tuple *(S,A,X,T)* where

- *S* and *A* are states and actions
- *X* is a a **parameter space**
- *T: S* × *A* × *X* → $\mathscr{D}$*S*

# Parametric Markov Decision Process <sub>pMDP</sub>



$$\mathbf{Pr}(\text{\textheart}) = \quad 0, \; \frac{2p - p^2}{2p^2 - 4p + 8}, \; -\frac{p^2}{2p - 4}, \; \frac{2p - p^2}{4p^2 - 8p + 8}, \; \frac{1}{8}\left(2\right.$$
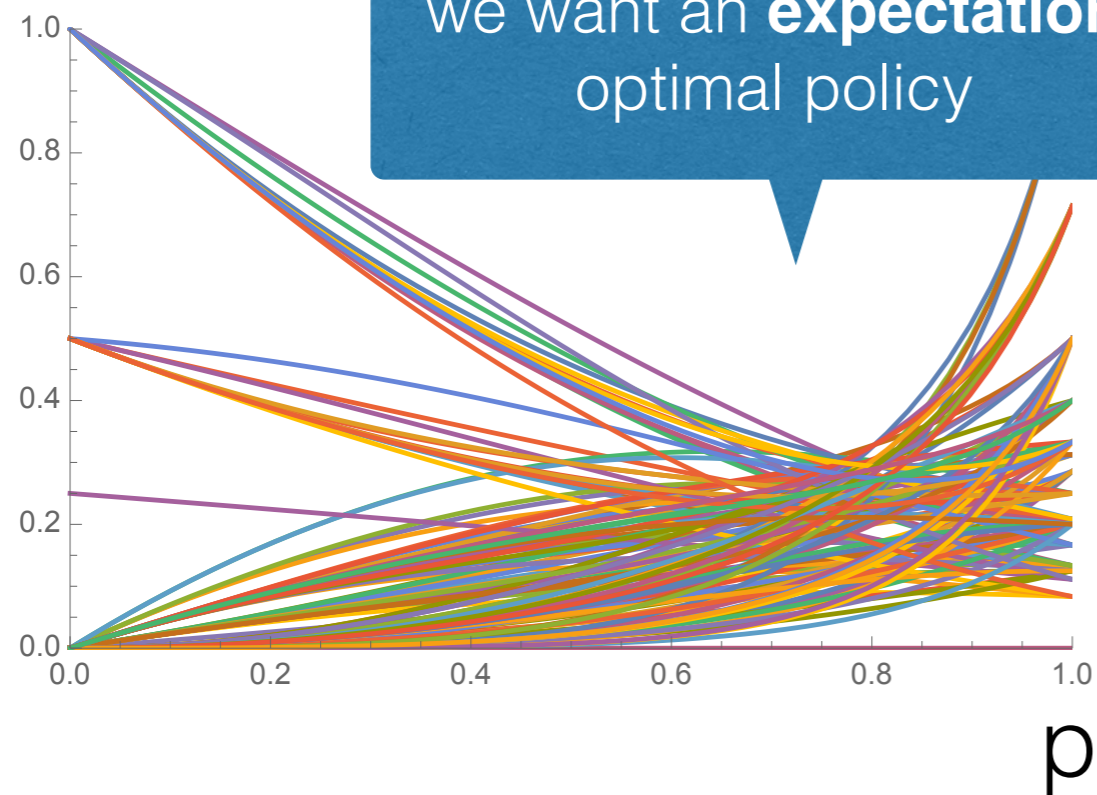
# Parametric Markov Decision Process <small>pMDP</small>



$$\mathbf{Pr}(\text{❤️}) = \quad 0, \; \frac{2p - p^2}{2p^2 - 4p + 8}, \; -\frac{p^2}{2p - 4}, \; \frac{2p - p^2}{4p^2 - 8p + 8}, \; \frac{1}{8}(2$$

# Expectation optimality



Pr(❤️)

we want an **expectation** optimal policy

p

Pr(❤️)

optimizes the area
(for uniform distribution on p)

p

30

# Expectation optimality

A policy $\pi$, initial distribution *i* and a distribution over the parameters *d*, gives us a  parametric probability space (Runs$_X$, Cones$_X$, **P**$_{\pi,i,d}$ ).


Runs$_X$    - disjoint union of the runs for all parameter values

Cones$_X$ - the $\sigma$ algebra generated by the disjoint union of cones for all parameter values

**P**$_{\pi,i,d}$       - the d-convex combination of the individual measures on cones

# Expectation optimal policy

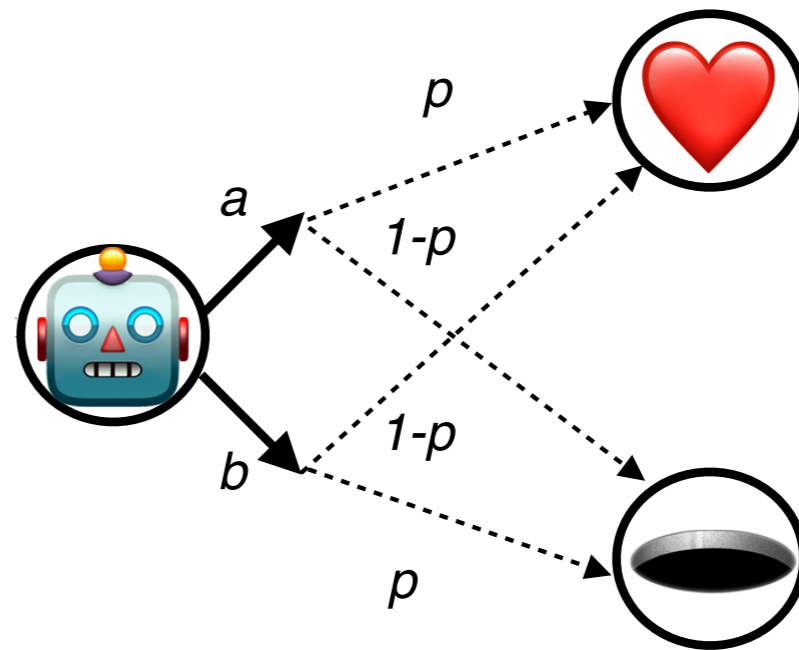An **expectation optimal policy** is a policy $\pi$ with

$$\mathbf{E}_{\pi,i,d}(r) = \sup_{\sigma} \mathbf{E}_{\sigma,i,d}(r)$$

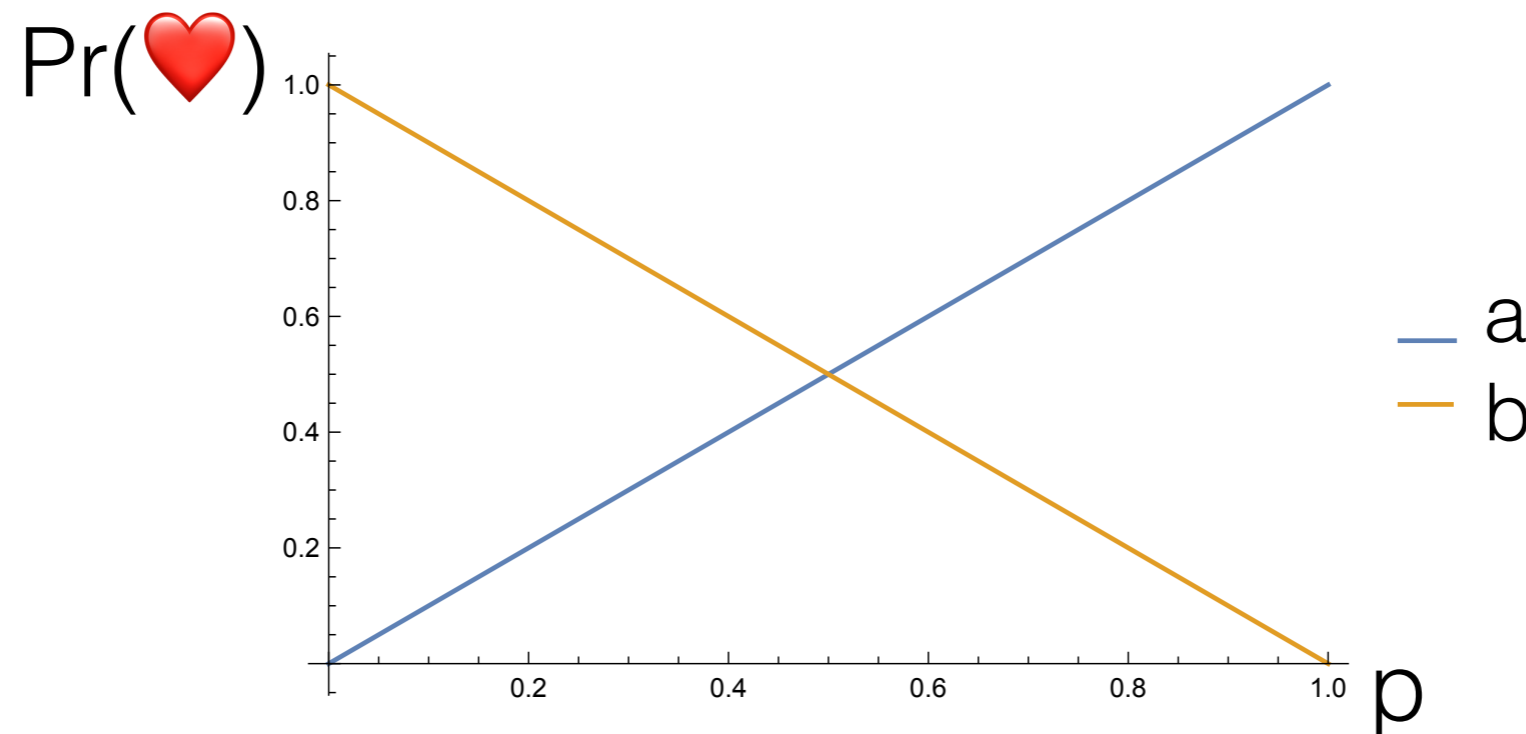An optimal policy not always exists, but **ε-optimal** do:

A policy is **expectation ε-optimal** if

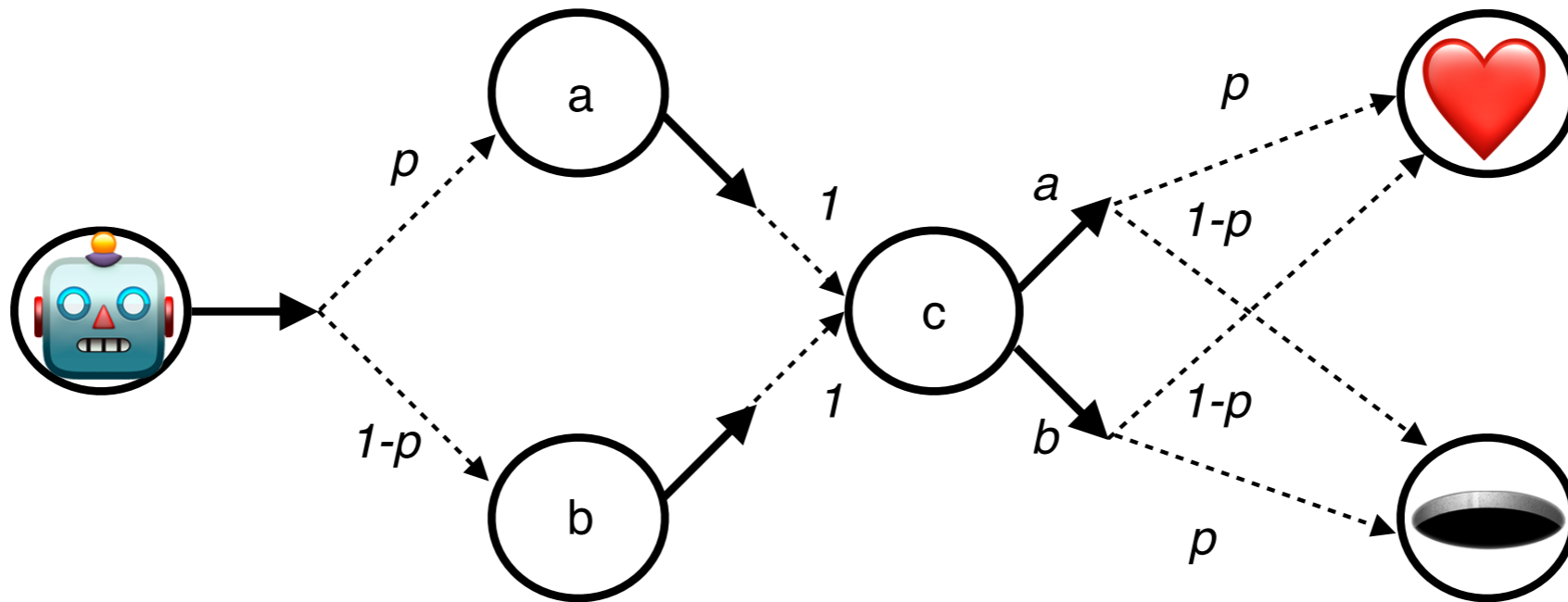$$\mathbf{E}_{\pi,i,d}(r) \text{ is ε-close to } \sup_{\sigma} \mathbf{E}_{\sigma,i,d}(r)$$
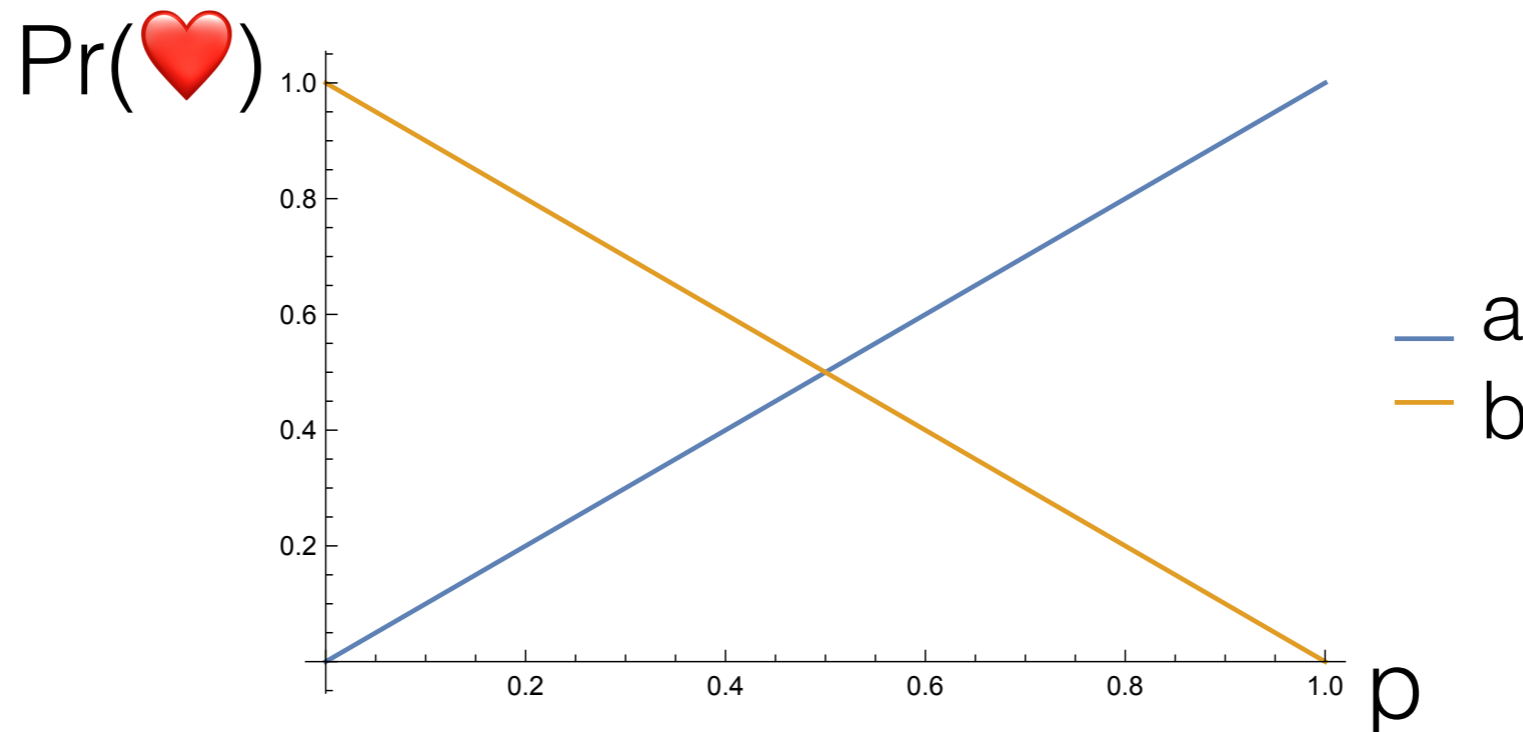
# Learner Example



**simple policies**

# Learner Example
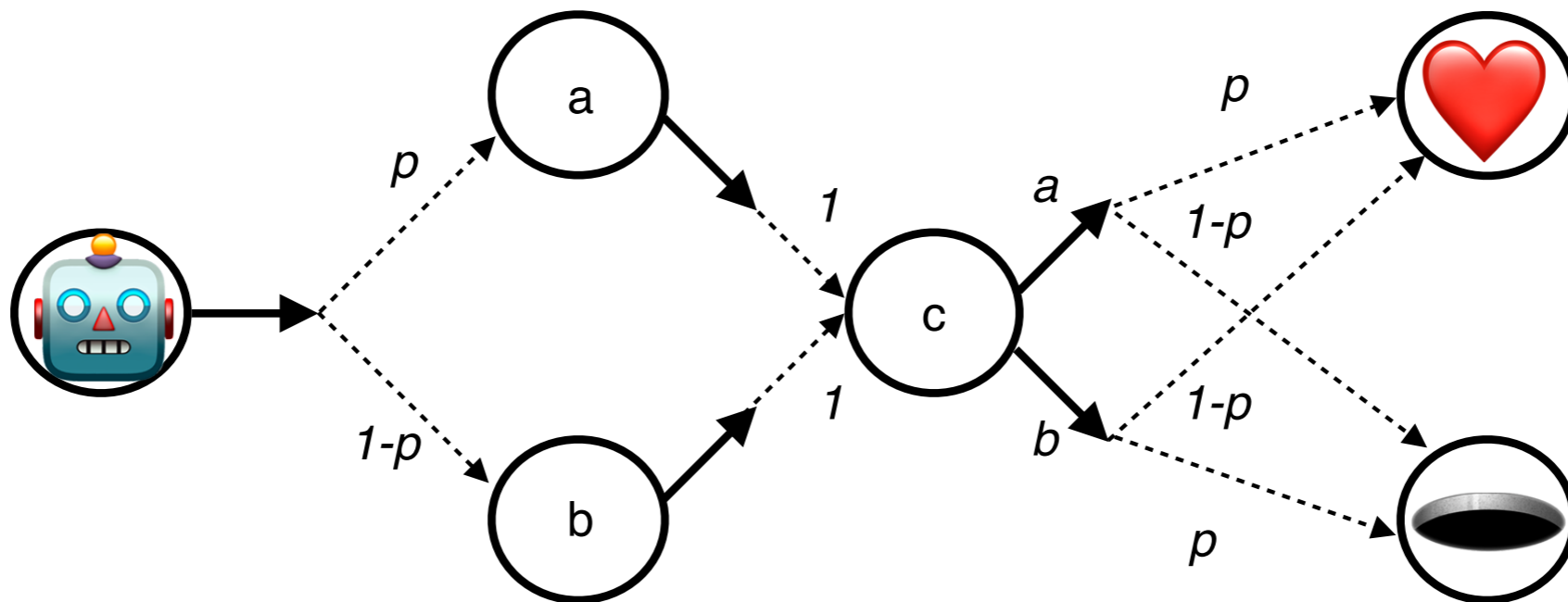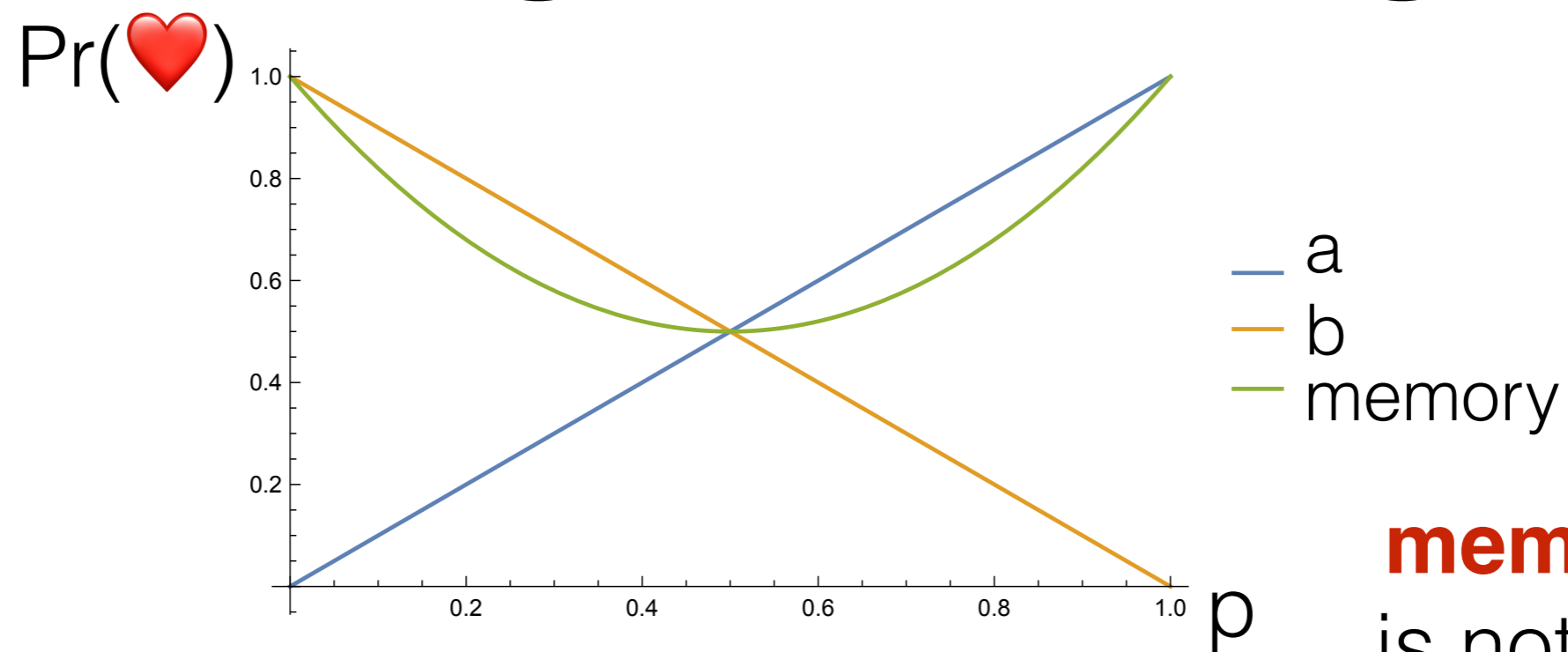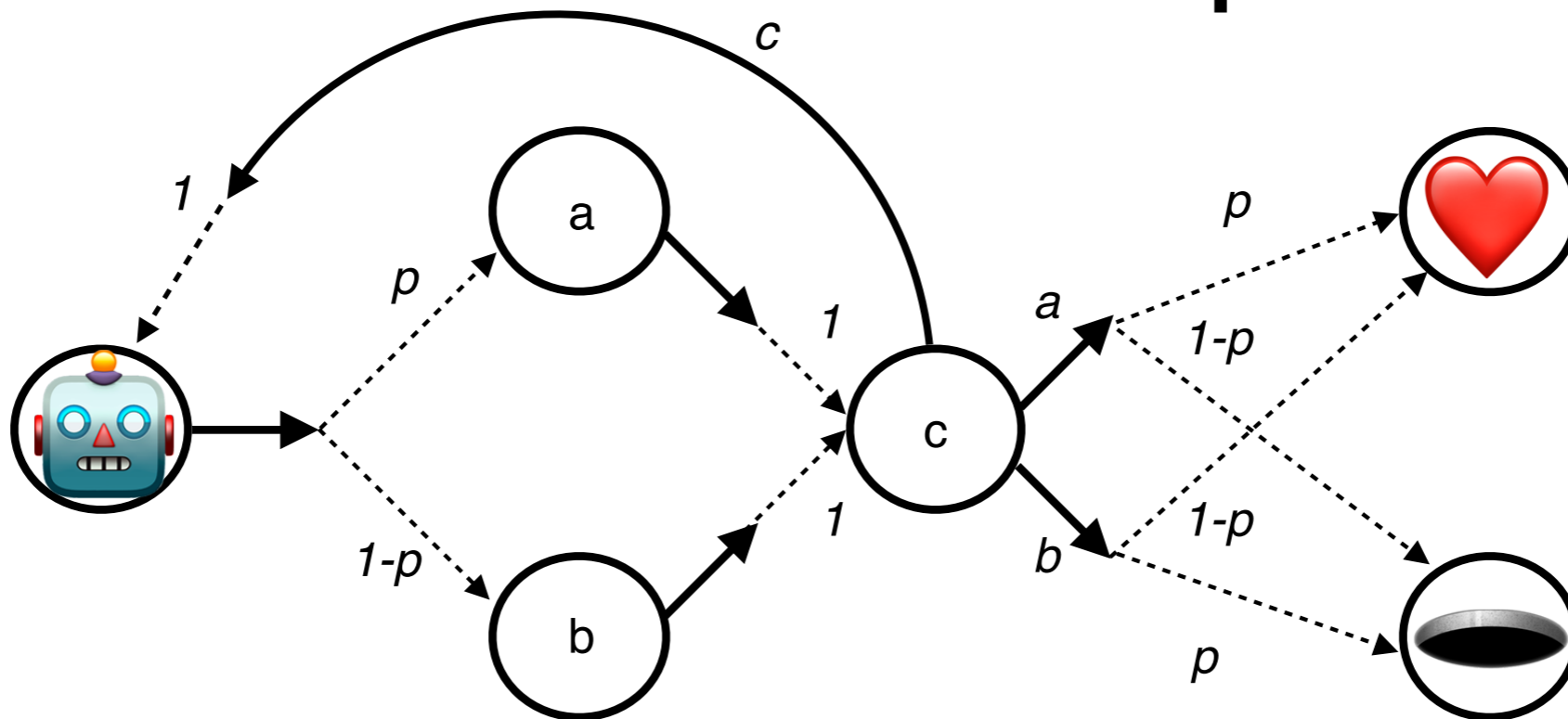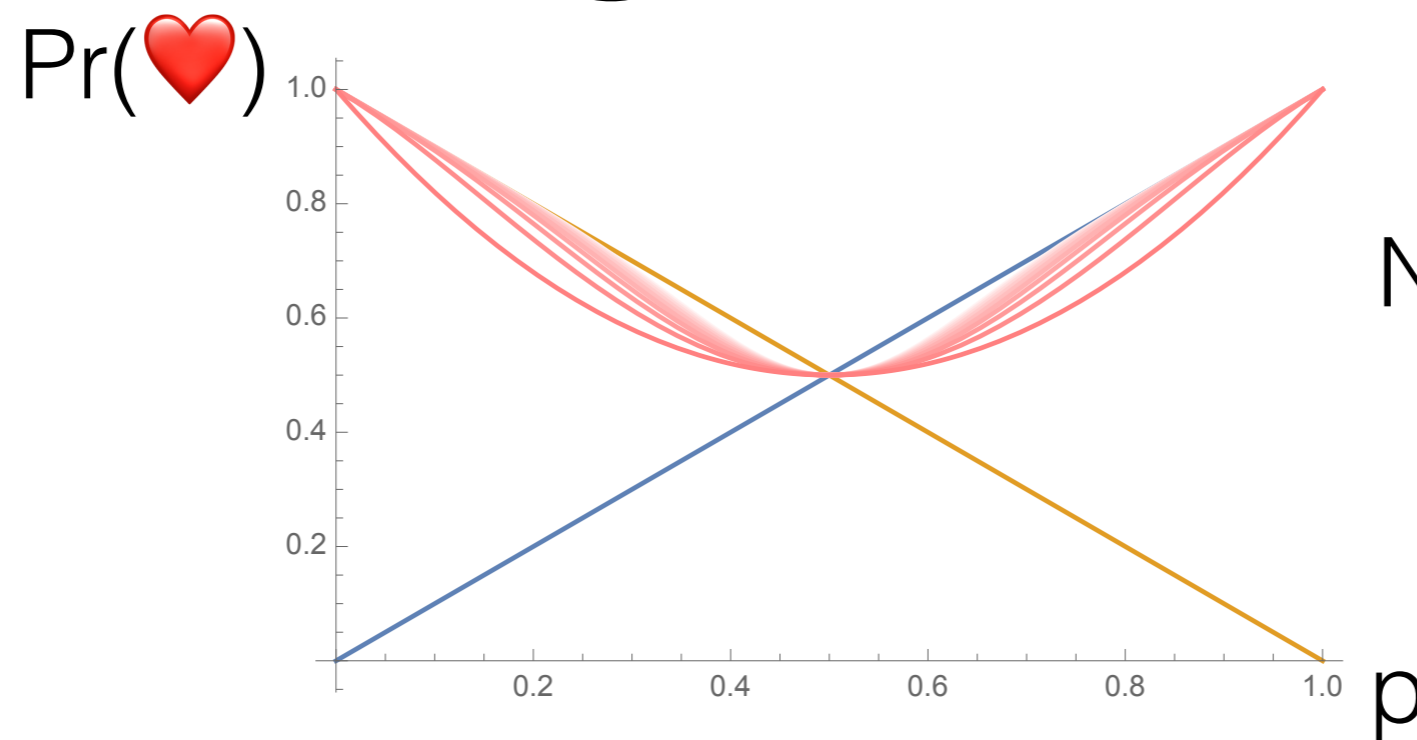


**simple policies**

# Learner Example



Pr(❤️)

**policies**

a
b
memory

**memoryless**
is not enough

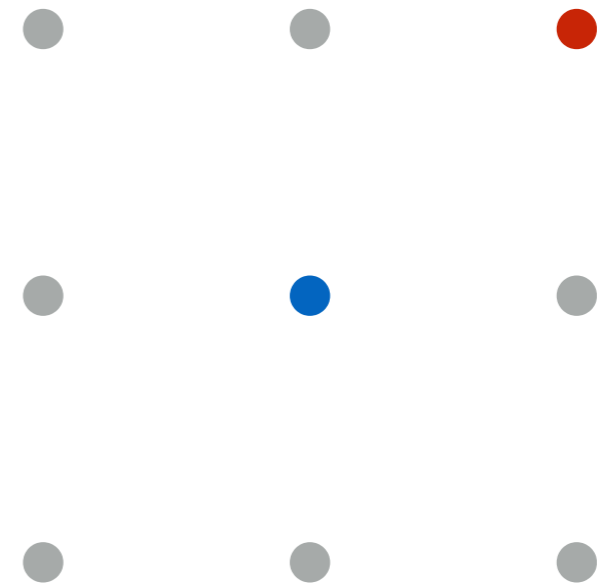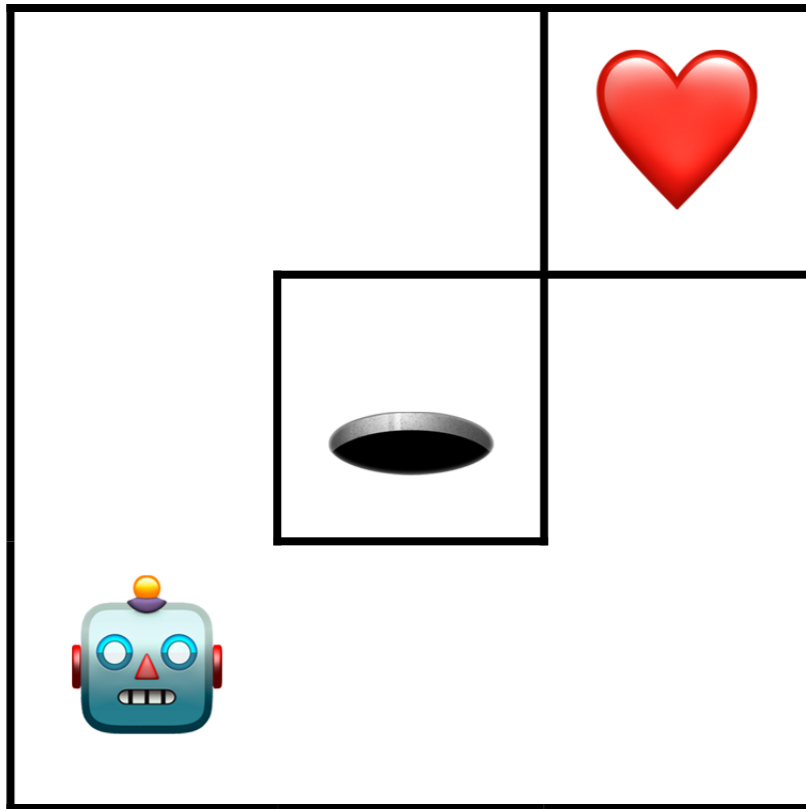# Learner Example



Pr(❤️)

**policies**

No **optimal** policy
but $\varepsilon$-**optimal**

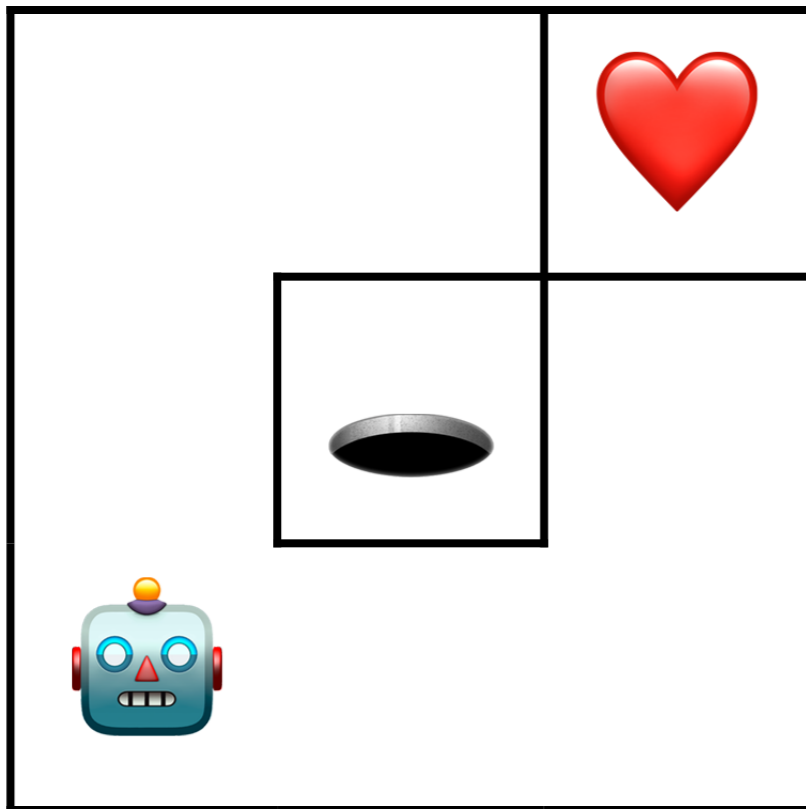**memoryless**
is not enough

# How to compute these policies?

# POMDP: Partially Observable MDP

# Partially Observable MDP POMDP

# Partially Observable MDP POMDP



A **POMDP** is a tuple $(S,A,T,\Omega,O)$ where

- $(S,A,T)$ is an MDP

- $\Omega$ is a set of **observations**
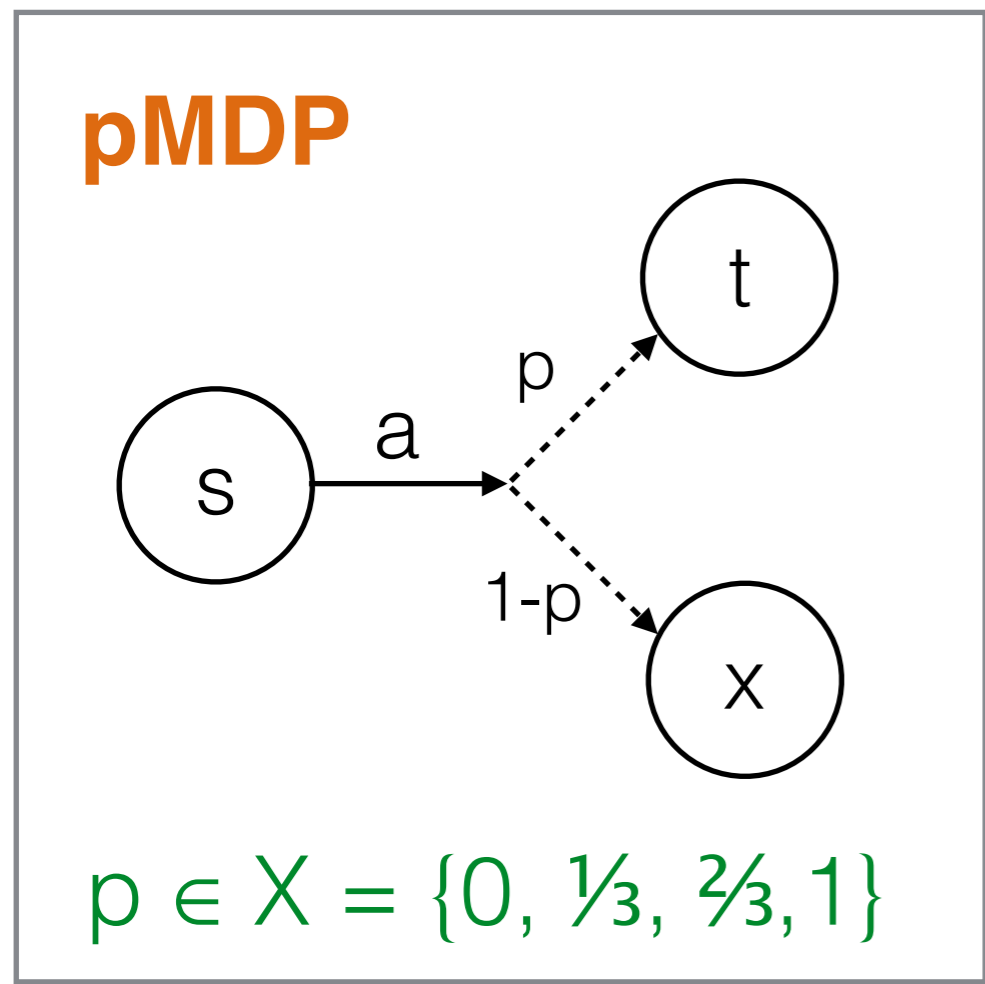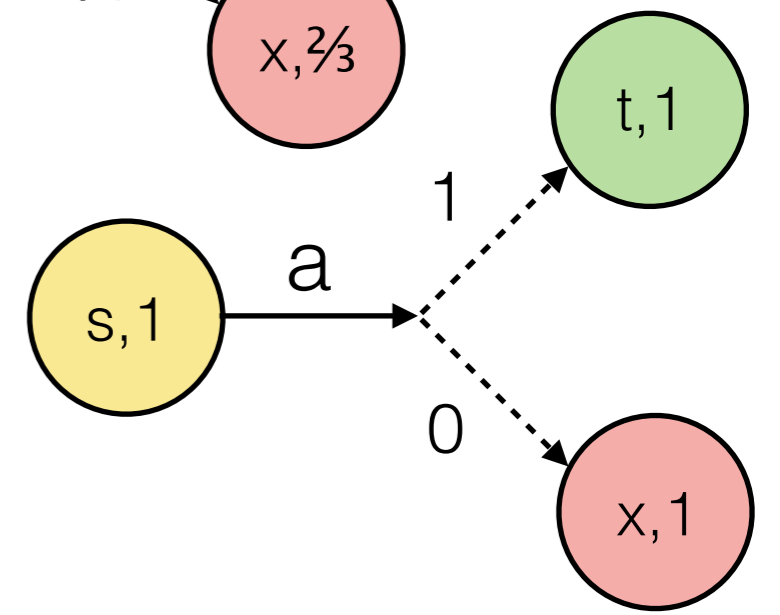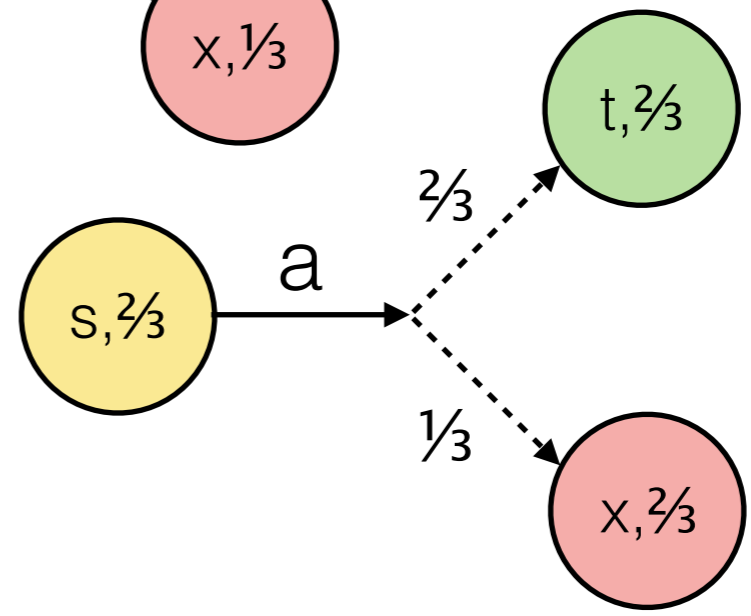
- $O: S \to \Omega$ is the **observation function**

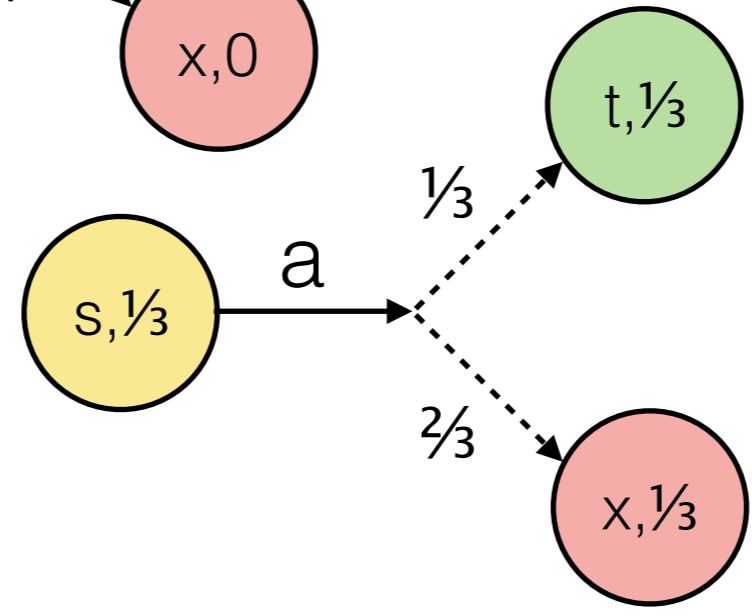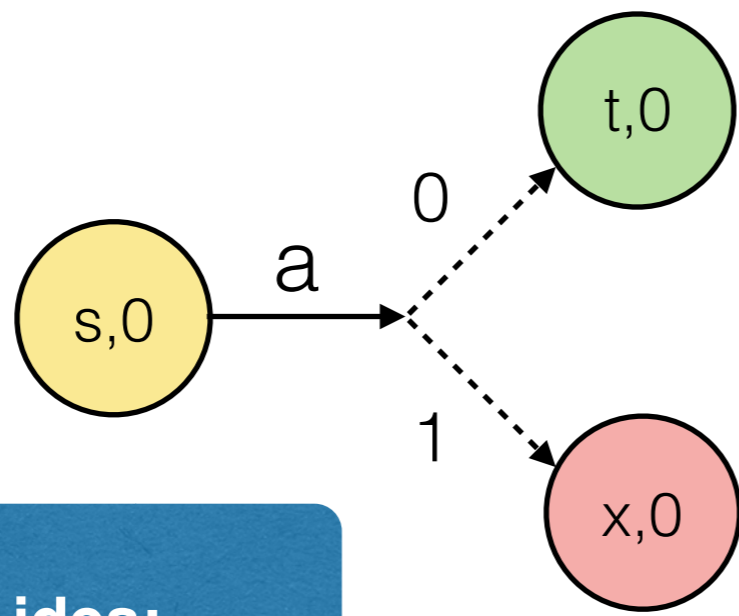A **POMDP policy** $\pi$ is a function

$\pi: (\Omega \times A)_* \times \Omega \to \mathcal{D}A$

**Encoding main idea:**
put parameter into POMDP states
observe only pMDP states

# POMDP

**Encoding main idea:**
put parameter into POMDP states
observe only pMDP states



## pMDP

$p \in X = \{0, \frac{1}{3}, \frac{2}{3}, 1\}$

# The Encoding

Given a **pMDP M** *(S,A,X,T)* we construct the **POMDP M'** *(S',A',T',Ω,O)*, where

$S' = S \times X$

$A' = A$

$T'((s,x),a)(s',x') = T(s,a)(x)(s') \cdot \delta_x(x')$

$\Omega = S$

$O((s,x)) = s$

**Note**: There is a 1-1 correspondence between the policies of **M** and **M'**.

# Correctness

Hence we can use off-the-shelf POMDP tools to compute expectation optimal pMDP policies.

**Theorem**:
Given a **pMDP M** and its **POMDP** encoding **M'**:

every **ε-optimal policy** of **M'** is an **ε-expectation optimal policy** for **M**, and vice versa.
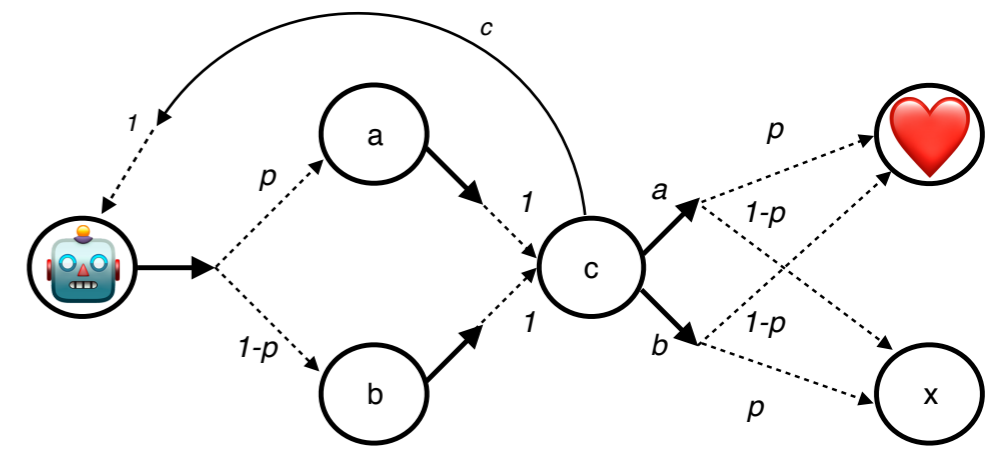
# Tools

- Work for finite horizon reward objectives

- Online and **Offline** algorithms

- **AI-Toolbox**

  - Incremental Pruning (IP)

  - Point Based Value Iteration (PBVI)

computing $\varepsilon$-optimal policies for infinite horizon POMDPs is undecidable

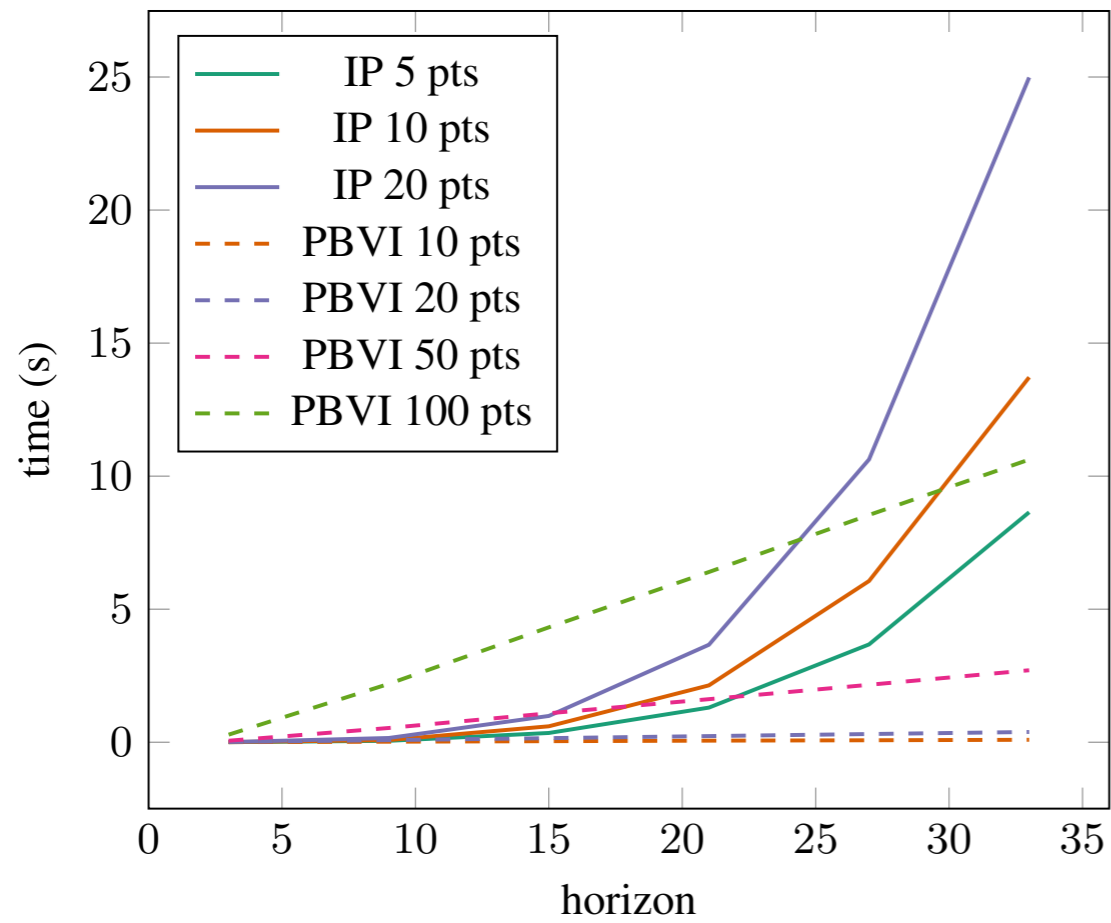Prism model → STORM → AI-Toolbox via Python interfaces
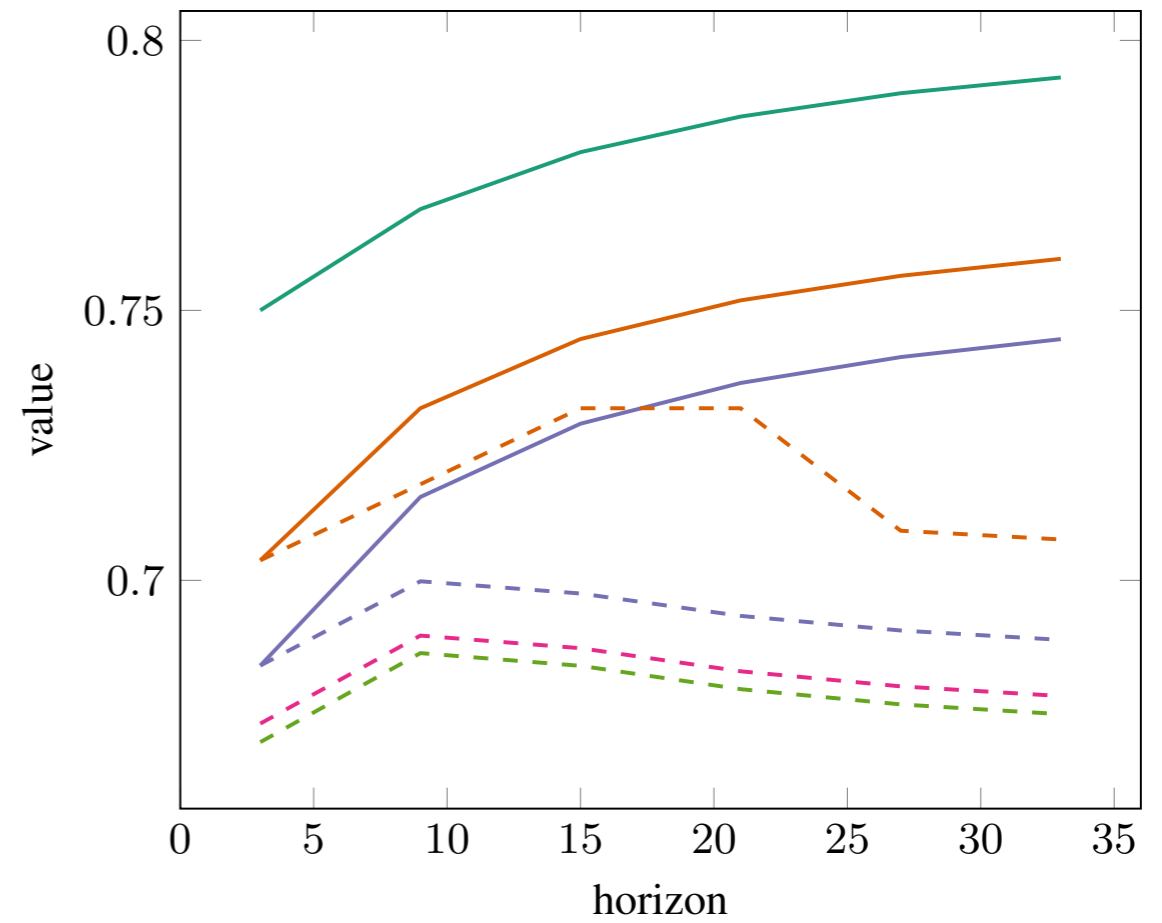
# Experimental Results

# Learner



uniform distribution over equidistant points in [0,1]
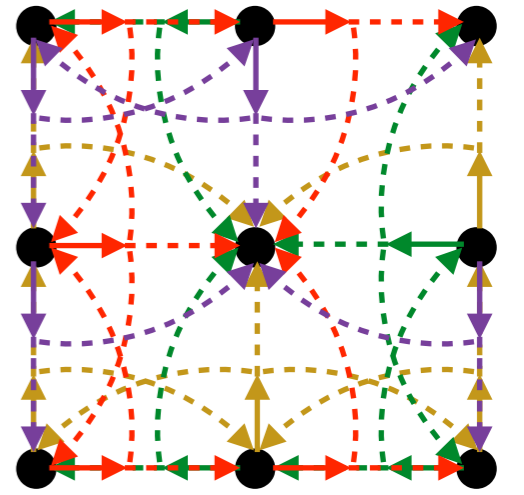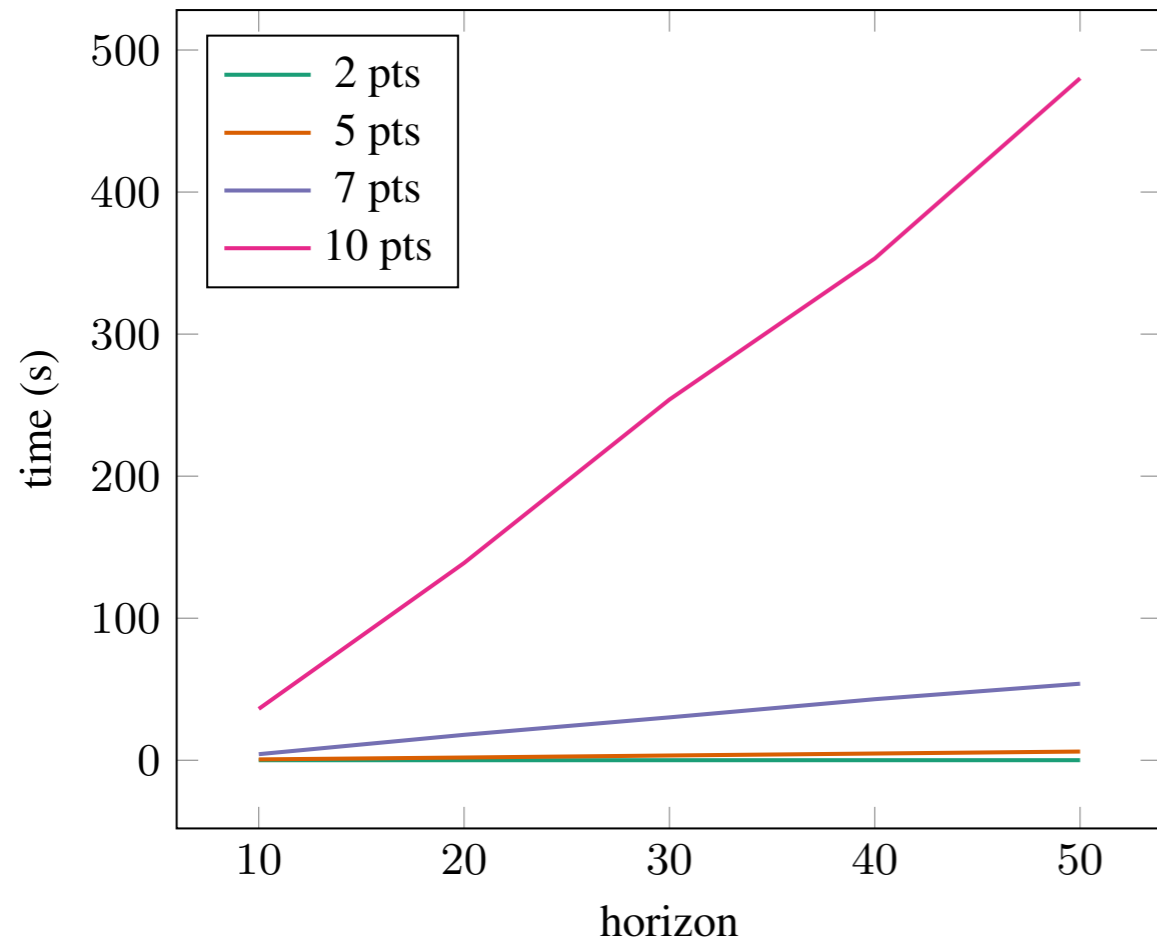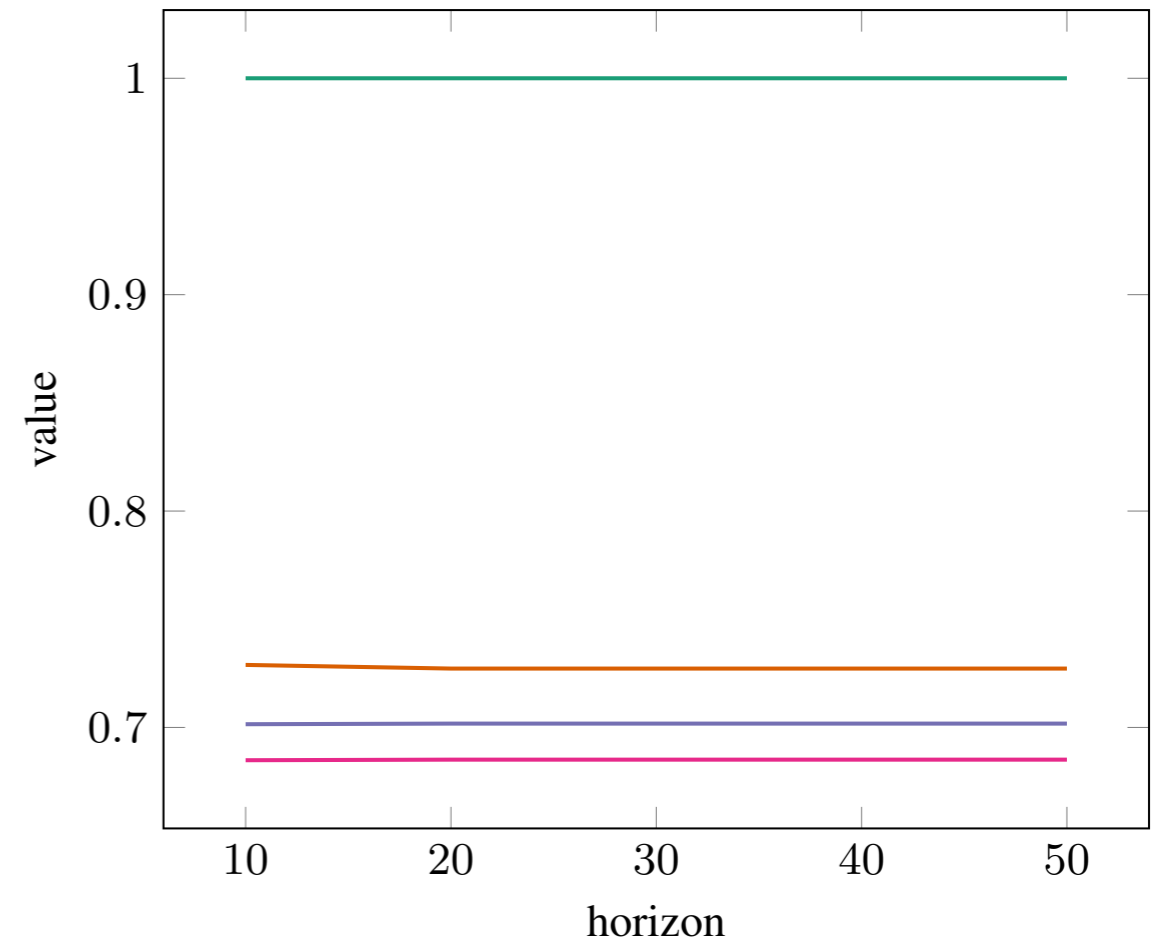


Runtime

Reachability

# Robot



Runtime



Reachability

# Summary

Finding **policies**

of a **parametric MDP**

that are **expectation optimal**

(over the **whole parameter space**)

amounts to solving a **suitable POMDP.**

We have a
proof of concept implementation

**github.com/sarming/pMDP-Toolbox**

**Thank You!**