

Model-Based Exploration in Monitored Markov Decision Processes



MSc defence seminar

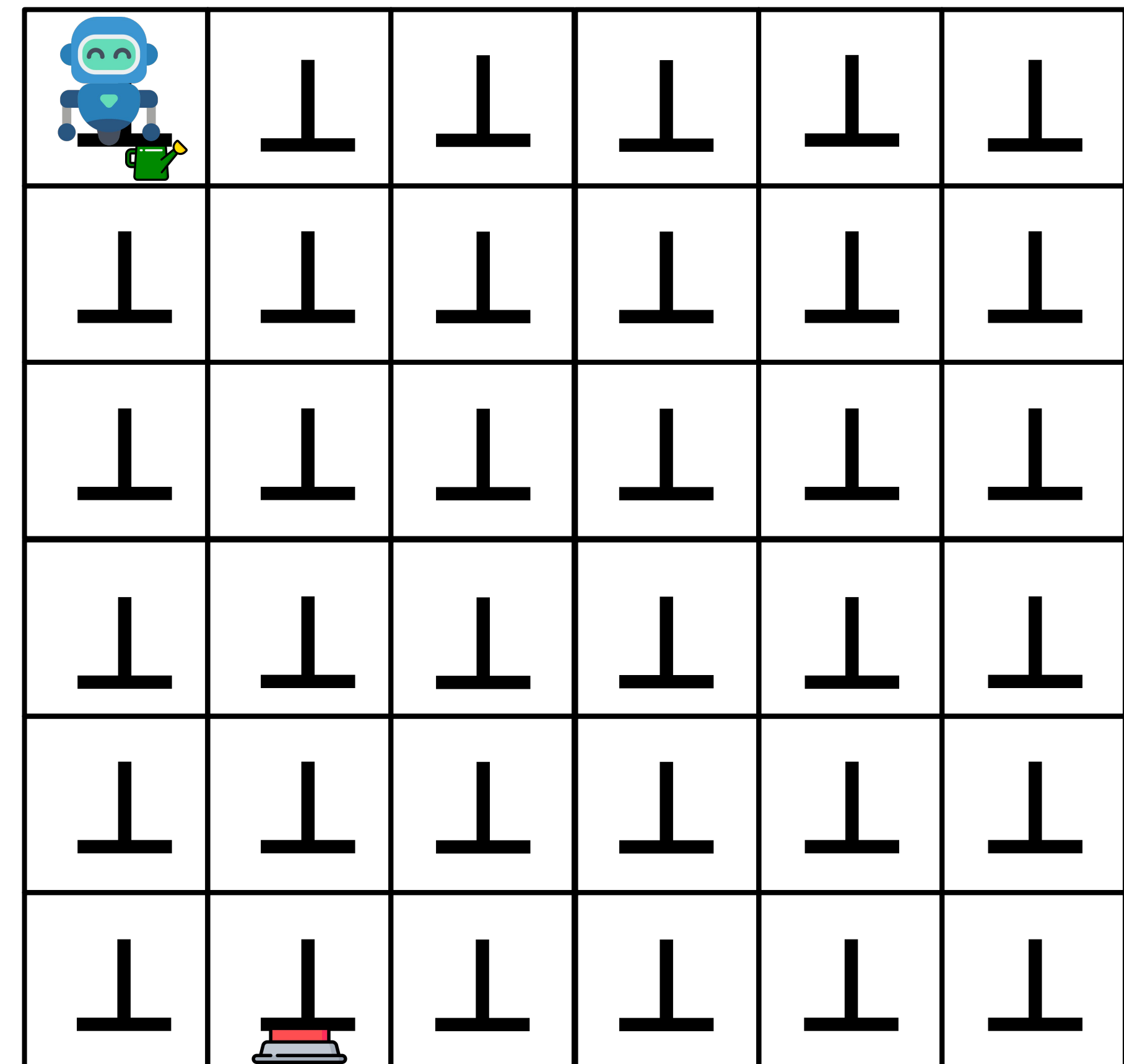
Alireza Kazemipour - August 21st, 2025

What should the agent do?



The agent's supposed to water .

The agent only knows there is at least one  that has the minimum reward a  that has the highest reward.





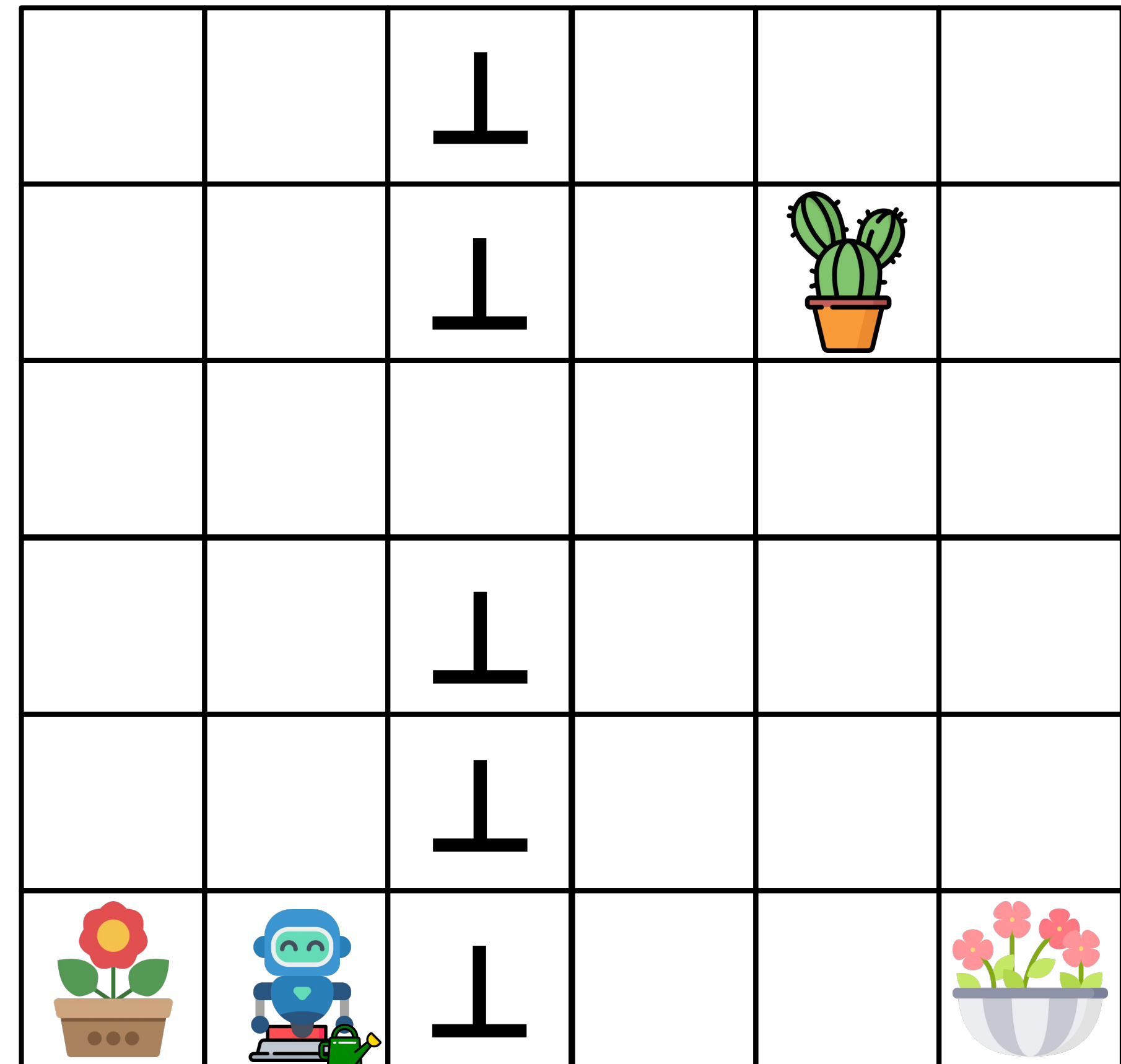
What should the agent do?

- 0.2





The agent's supposed to water .

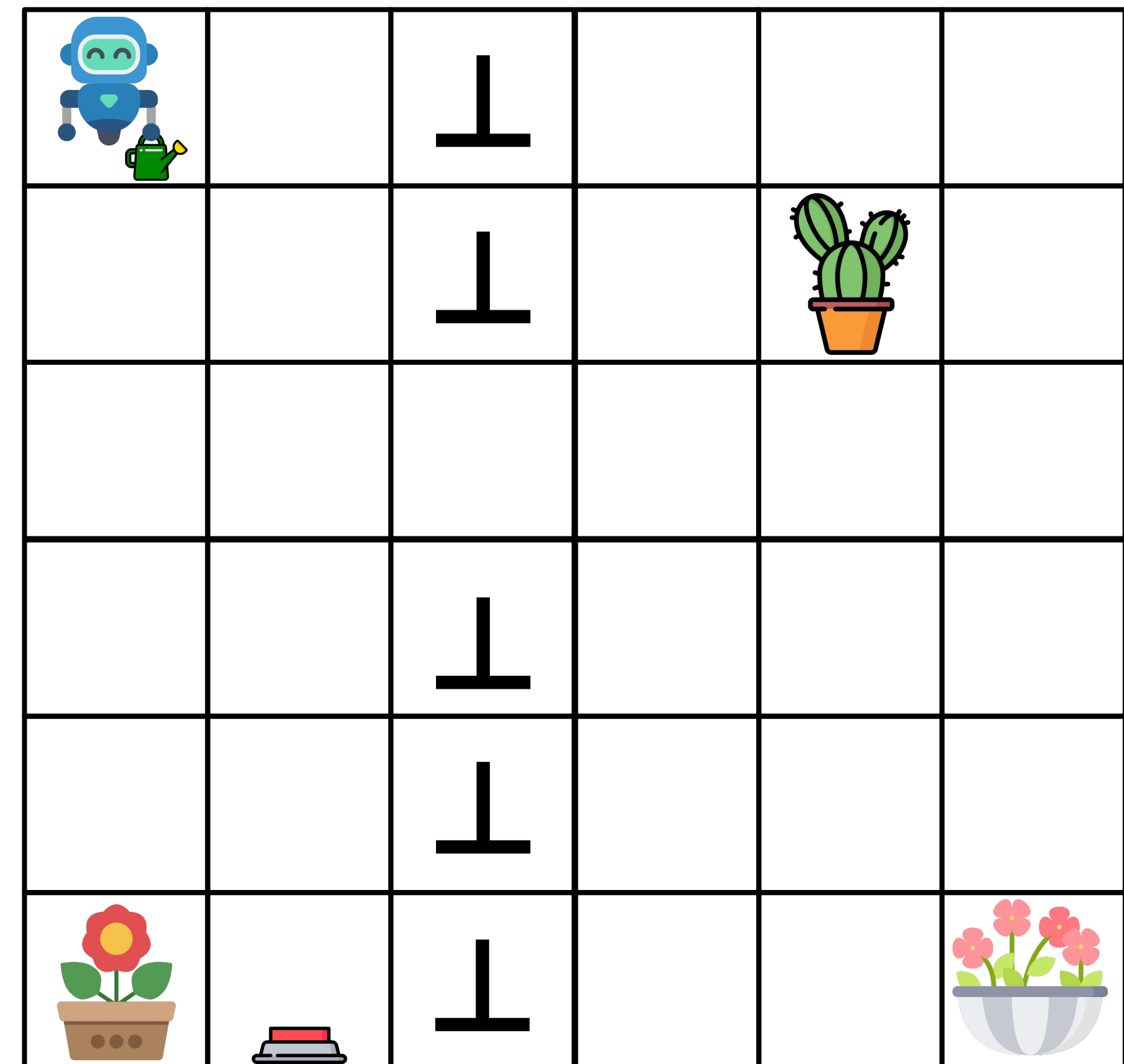
The agent only knows there is at least one that has the minimum reward a  that has the highest reward. 



What should the agent do?



The agent's supposed to water .

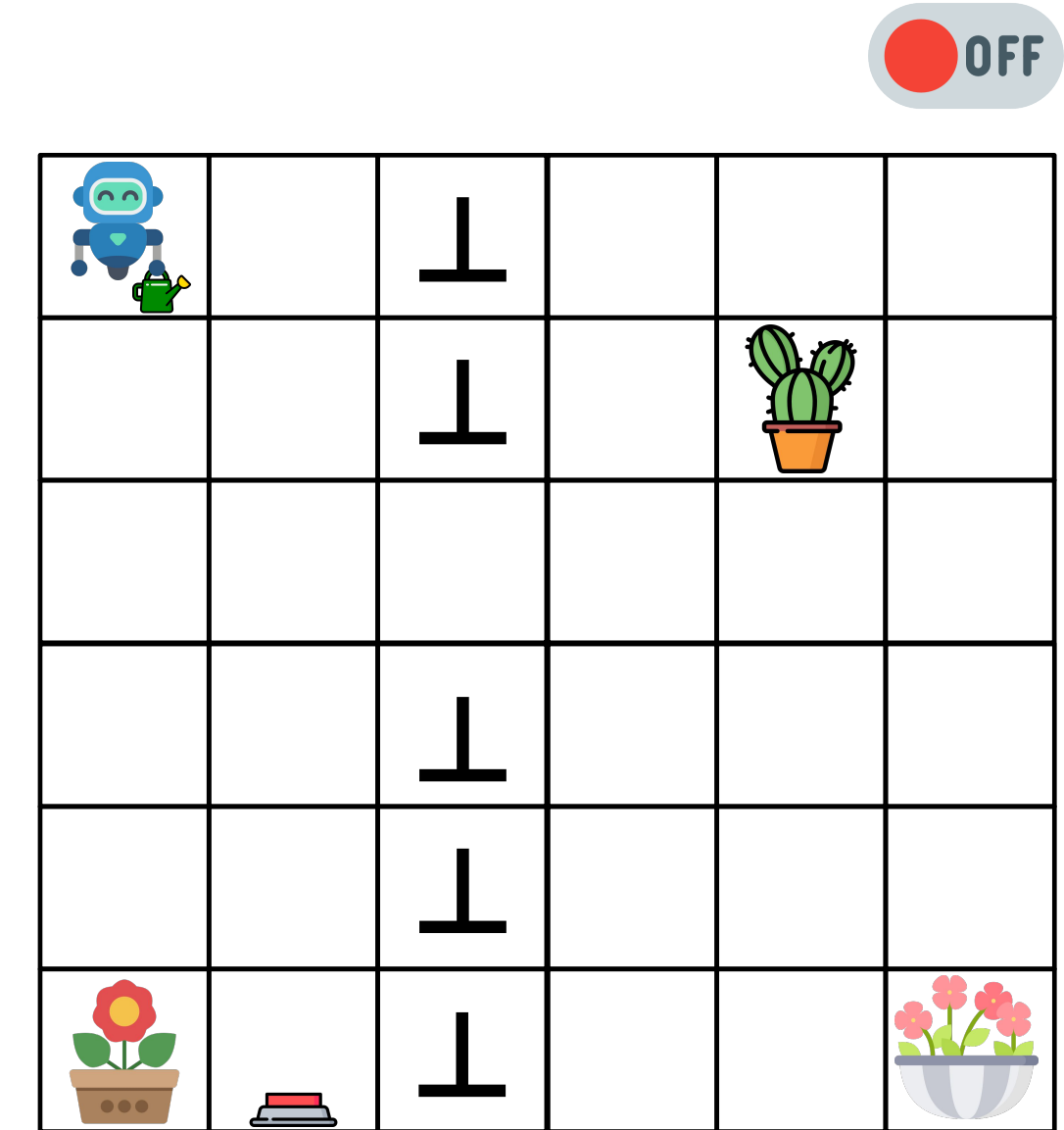
The agent only knows there is at least one  that has the minimum reward a  that has the highest reward.



What should the agent do?

The agent's supposed to water .

The agent only knows there is at least one  that has the minimum reward a  that has the highest reward.



I'm going to answer the question in the slide's title in this talk.

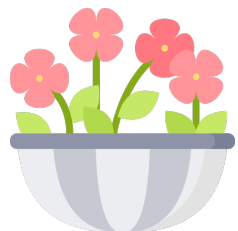
What should the agent do?

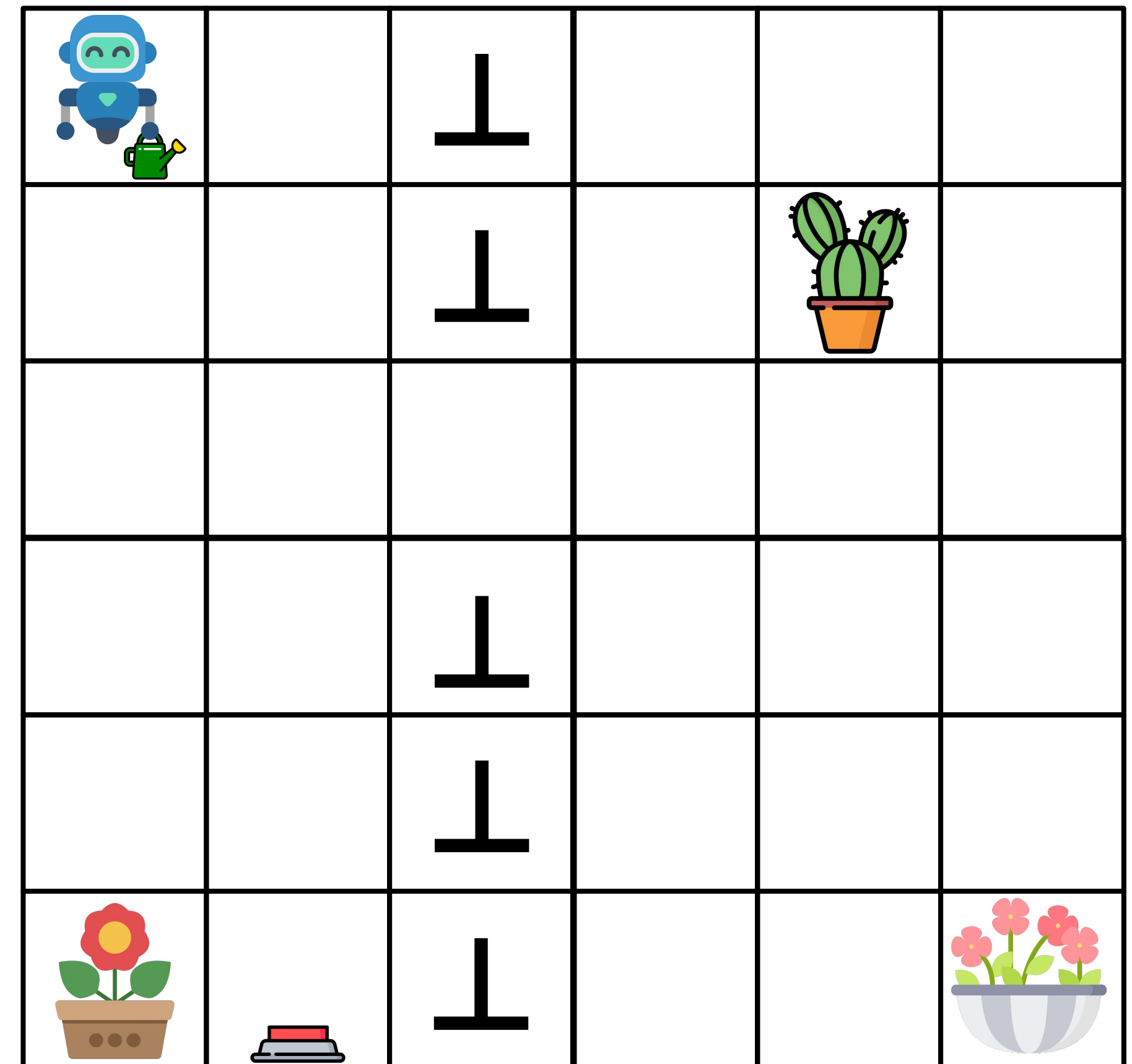
Breaking the overarching question into subproblems



1. How to detect \perp cells from all the others?

2. How to deal with \perp cells?

3. Can the agent be efficient in watering  while not impacting (1) and (2)?



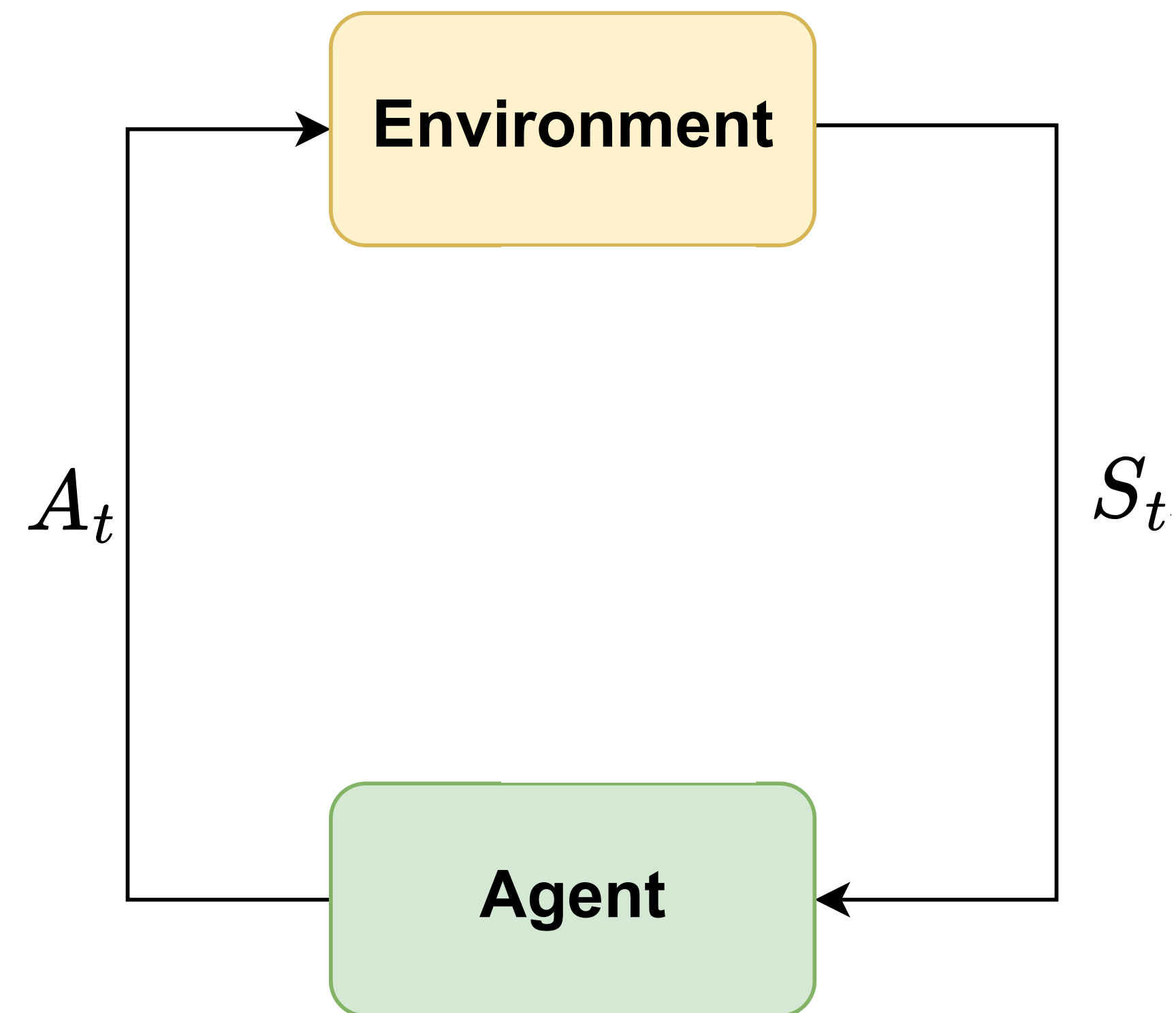
This talk

- Review:
 - Markov Decision Processes
 - Model-Based Interval Estimation with Exploration Bonus (MBIE-EB)
- Problem setting:
 - Monitored Markov Decision Processes
- My proposed solution: Monitored MBIE-EB
 - Theoretical performance
 - Empirical performance
- List of contributions
- Future work
- Acknowledgement

Review

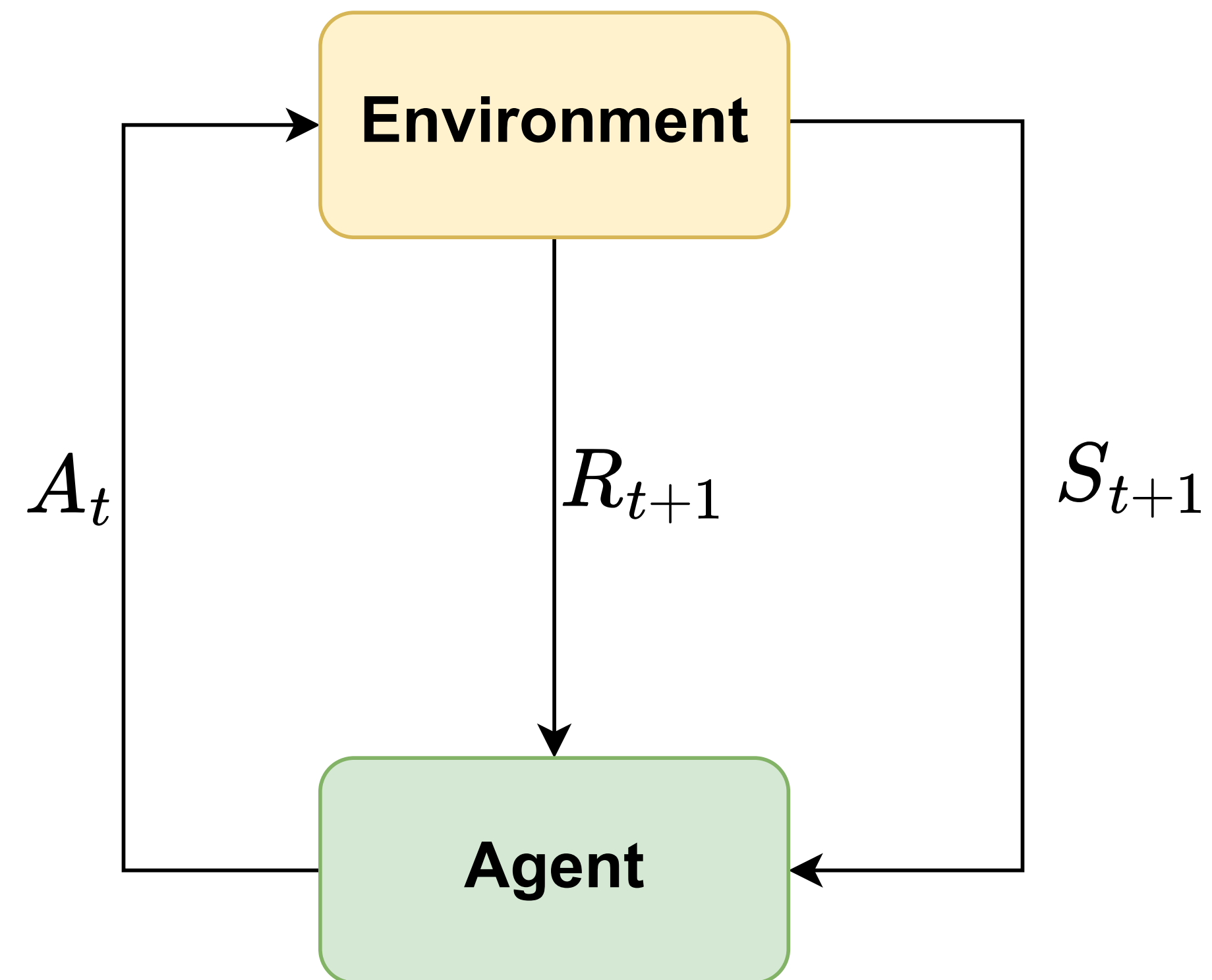
Markov decision processes (MDPs)

A typical mathematical model of interaction in RL



Markov decision processes (MDPs)

A typical mathematical model of interaction in RL



Markov decision processes (MDPs)

A typical mathematical model of interaction in RL

$$\text{The goal: } \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \right]$$

A finite MDP: $\langle \mathcal{S}, \mathcal{A}, r, p, \gamma \rangle$

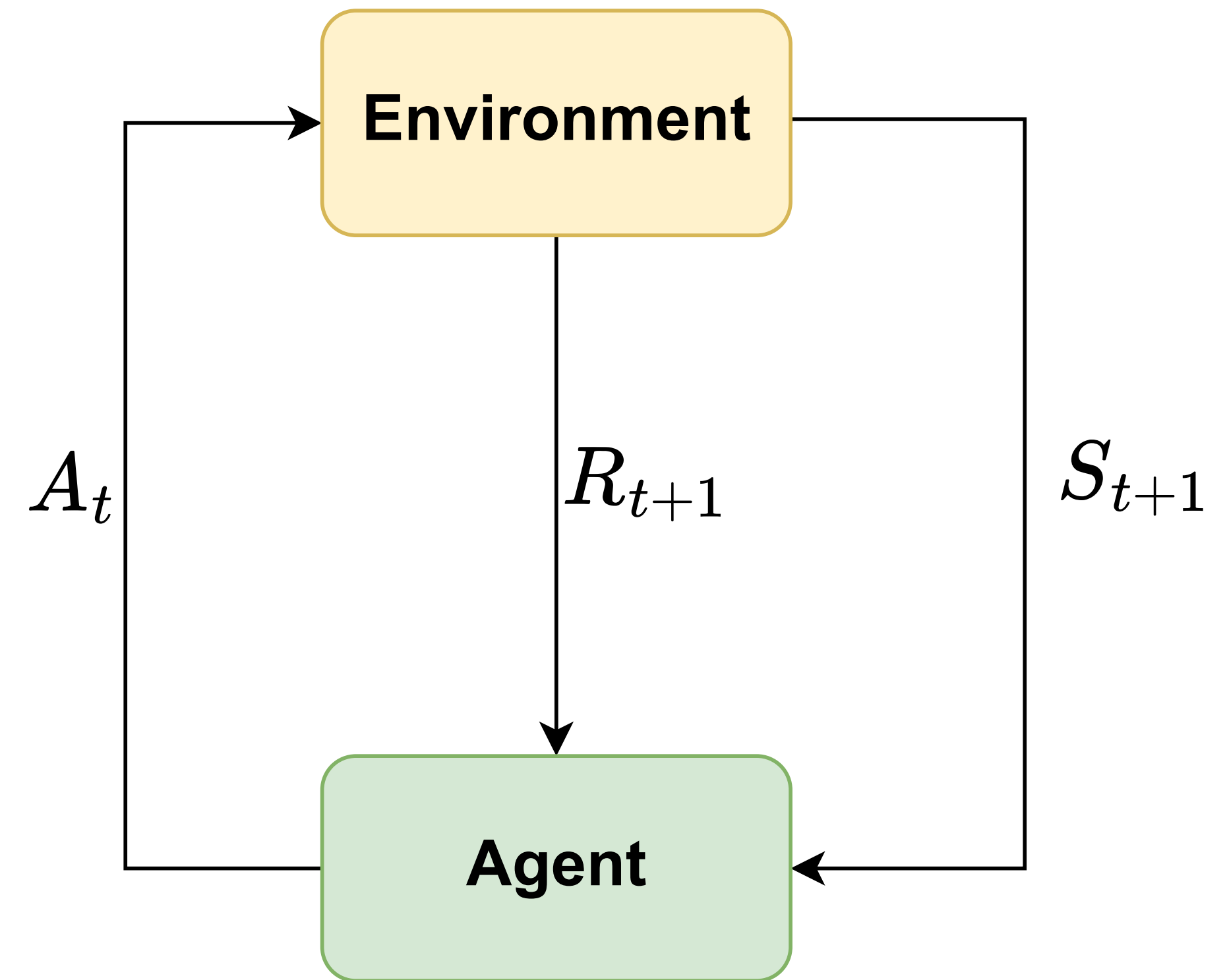
\mathcal{S} is the state space

\mathcal{A} is the action space

r is the expected immediate reward

p is the transition dynamics

$$0 \leq \gamma < 1$$



Markov decision processes (MDPs)

How to maximize the expected discounted return using models?

Follow

$$Q^*(s, a) = r(s, a) + \gamma \sum_{s'} p(s' | s, a) V^*(s').$$

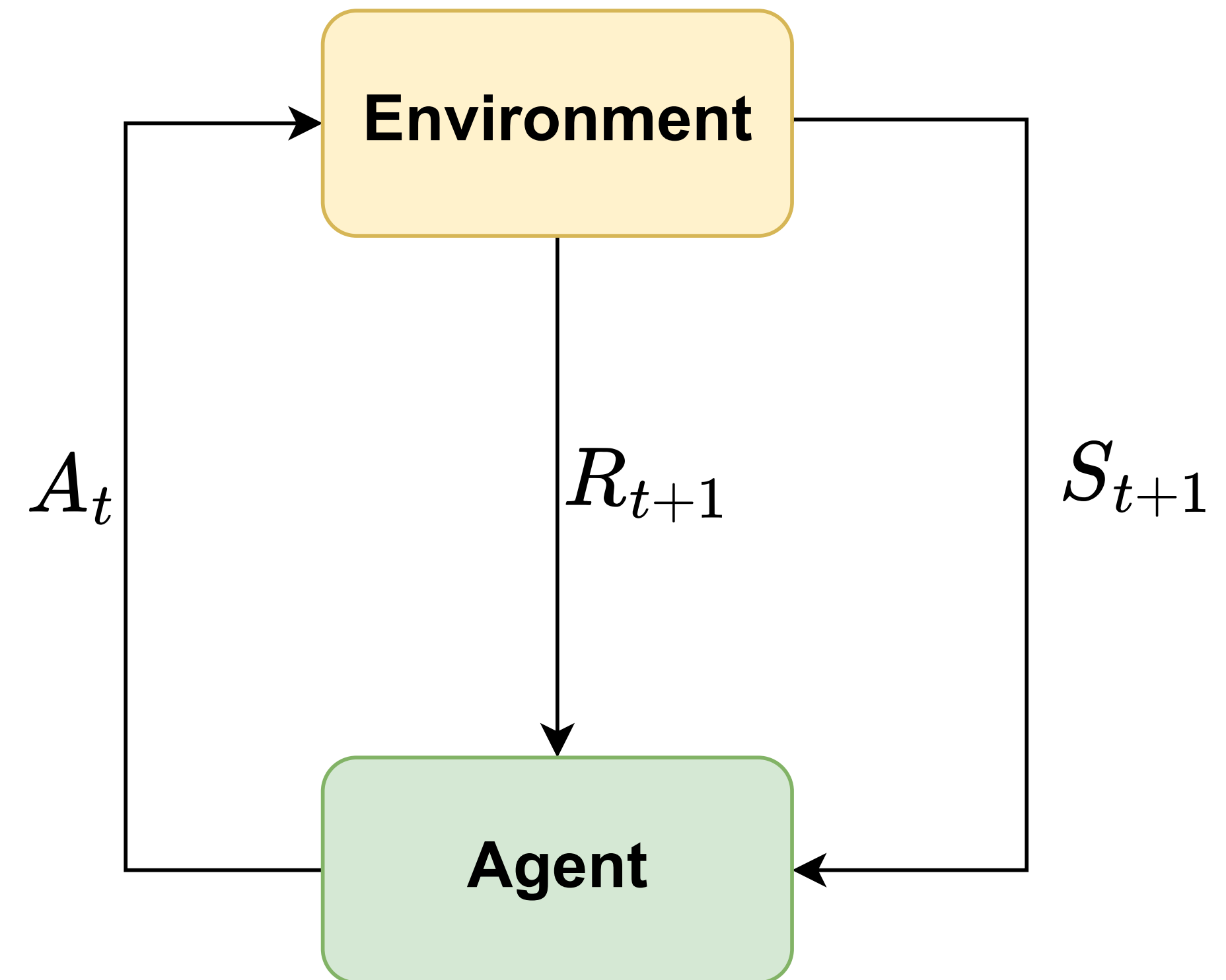
The model-based learning's challenge:

**We know sample estimates \hat{r} , and \hat{p} ,
but we don't know the true r and p ! 😭**

One solution:

Using measures on how uncertain we are about \hat{r} , and \hat{p} .

If we are confident about the quality of sample estimates, then we are golden.



Measuring uncertainty

- Suppose you have n samples. Then

$$\text{distance}(\text{empirical mean, true mean}) \leq \frac{\beta}{\sqrt{n}}$$

- If you particularly have n Bernoulli samples. Then

$$\text{distance}(\text{empirical mean, true mean}) \leq \frac{\beta}{n}$$

for sufficiently large value of β .

$\frac{\beta}{\sqrt{n}}$, and $\frac{\beta}{n}$ measure the
uncertainty.

Let \hat{Q} denote the action-value functions we get using \hat{r} , and \hat{p} .

Model-based interval estimation with exploration bonus (MBIE-EB)

There is an algorithm called **MBIE-EB**¹ that is greedy with respect to:

$$\hat{Q}(s, a)$$

1. A. Strehl, et al. “An analysis of model-based Interval Estimation for Markov Decision Processes,”
(Journal of Computer and System Sciences '08)

Model-based interval estimation with exploration bonus (MBIE-EB)

There is an algorithm called **MBIE-EB**¹ that is greedy with respect to:

$$\hat{Q}(s, a) + \underbrace{\frac{\beta_1}{\sqrt{n}}}_{\text{uncertainty of } \hat{r}}$$

1. A. Strehl, et al. “An analysis of model-based Interval Estimation for Markov Decision Processes,”
(Journal of Computer and System Sciences '08)

Model-based interval estimation with exploration bonus (MBIE-EB)

There is an algorithm called **MBIE-EB**¹ that is greedy with respect to:

$$\hat{Q}(s, a) + \underbrace{\frac{\beta_1}{\sqrt{n}}}_{\text{uncertainty of } \hat{r}} + \underbrace{\frac{\beta_2}{\sqrt{n}}}_{\text{uncertainty of } \hat{p}}$$

number visits to (s, a)

1. A. Strehl, et al. “An analysis of model-based Interval Estimation for Markov Decision Processes,” (Journal of Computer and System Sciences ’08)

Model-Based interval estimation with exploration bonus (MBIE-EB)

There is an algorithm called **MBIE-EB** that is greedy with respect to:

$$\hat{Q}(s, a) + \underbrace{\frac{\beta_1}{\sqrt{n}}}_{\text{uncertainty of } \hat{r}} + \underbrace{\frac{\beta_2}{\sqrt{n}}}_{\text{uncertainty of } \hat{p}}$$

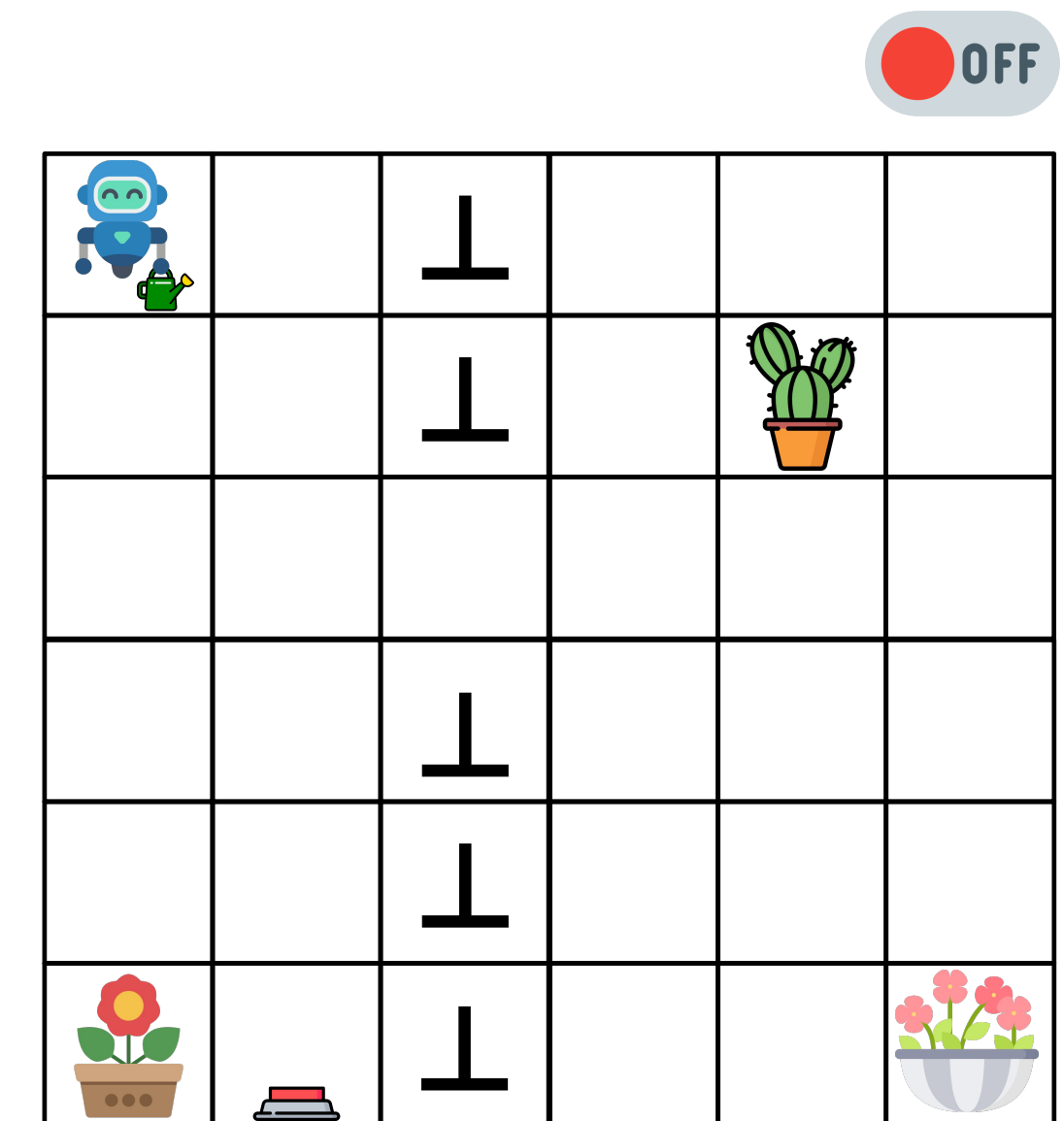
number visits to (s, a)

MBIE-EB is also efficient since it finds an ϵ -optimal policy in following number of time steps:

$$\tilde{O}\left(\frac{|\mathcal{S}|^2 |\mathcal{A}|}{\epsilon^3 (1 - \gamma)^6}\right)$$

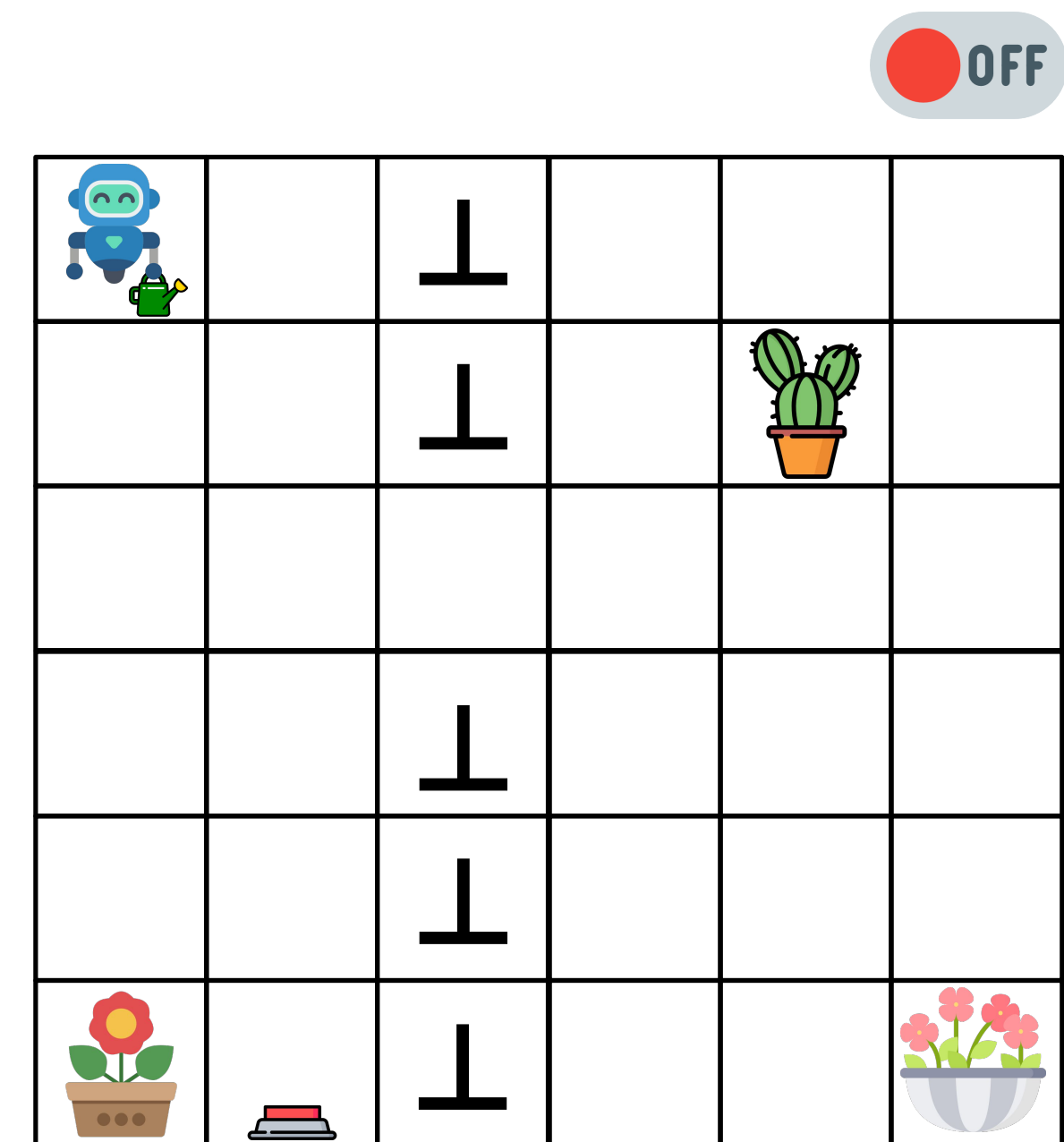
Takeaways so far:

1. A good algorithm like MBIE-EB uses **bonuses** as measures of **uncertainty**
2. We're interested in solving



Problem setting

MDPs cannot model



But Monitored MDPs¹ can!

1. S. Parisi, et al. “Monitored Markov Decision Processes,” (AAMAS '24)

Monitored Markov decision processes (Mon-MDPs)

An extension of MDPs to cover partial observability of rewards

```
graph TD; Environment[Environment] --- Monitor[Monitor]; Monitor --- Agent[Agent];
```

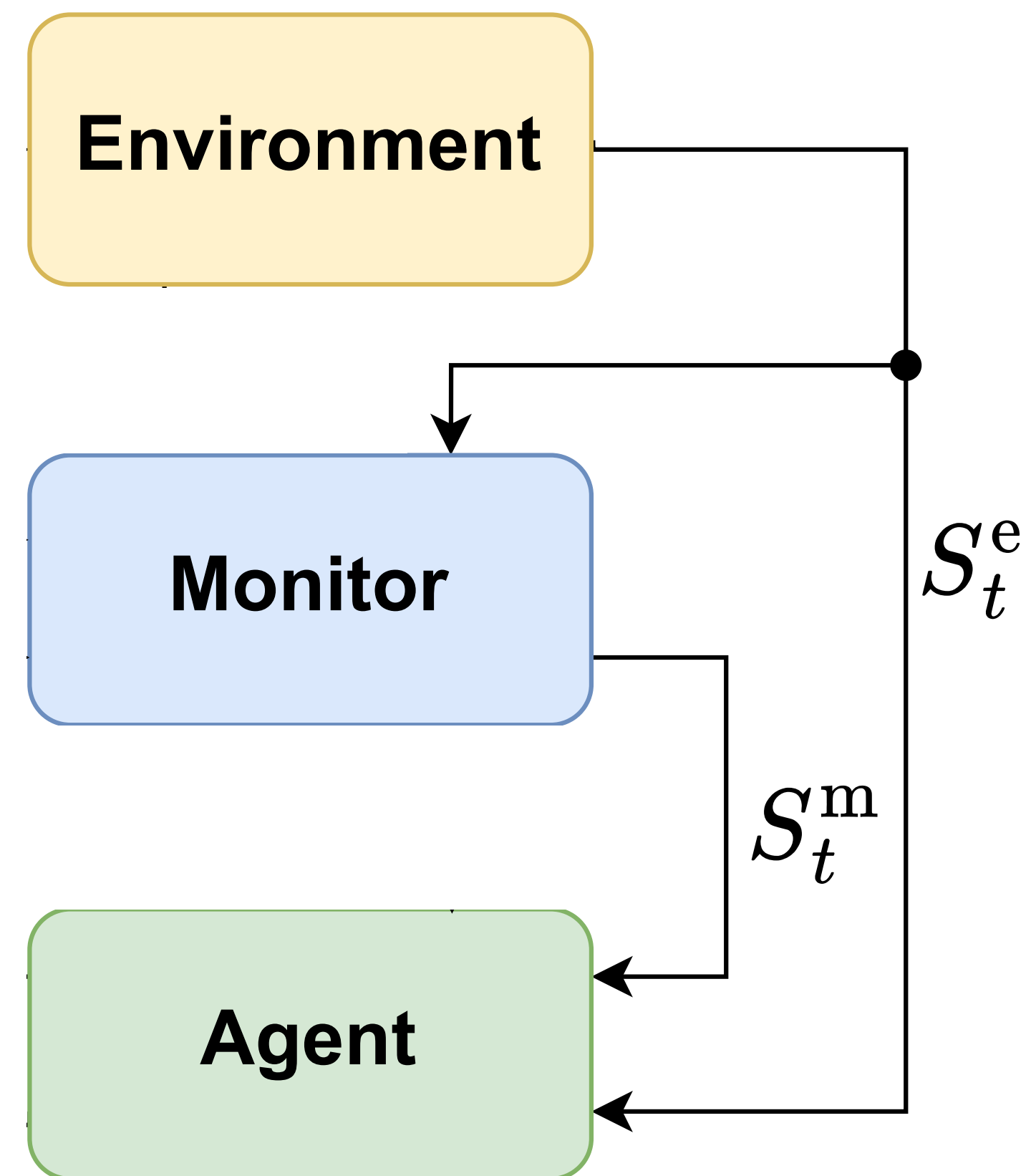
Environment

Monitor

Agent

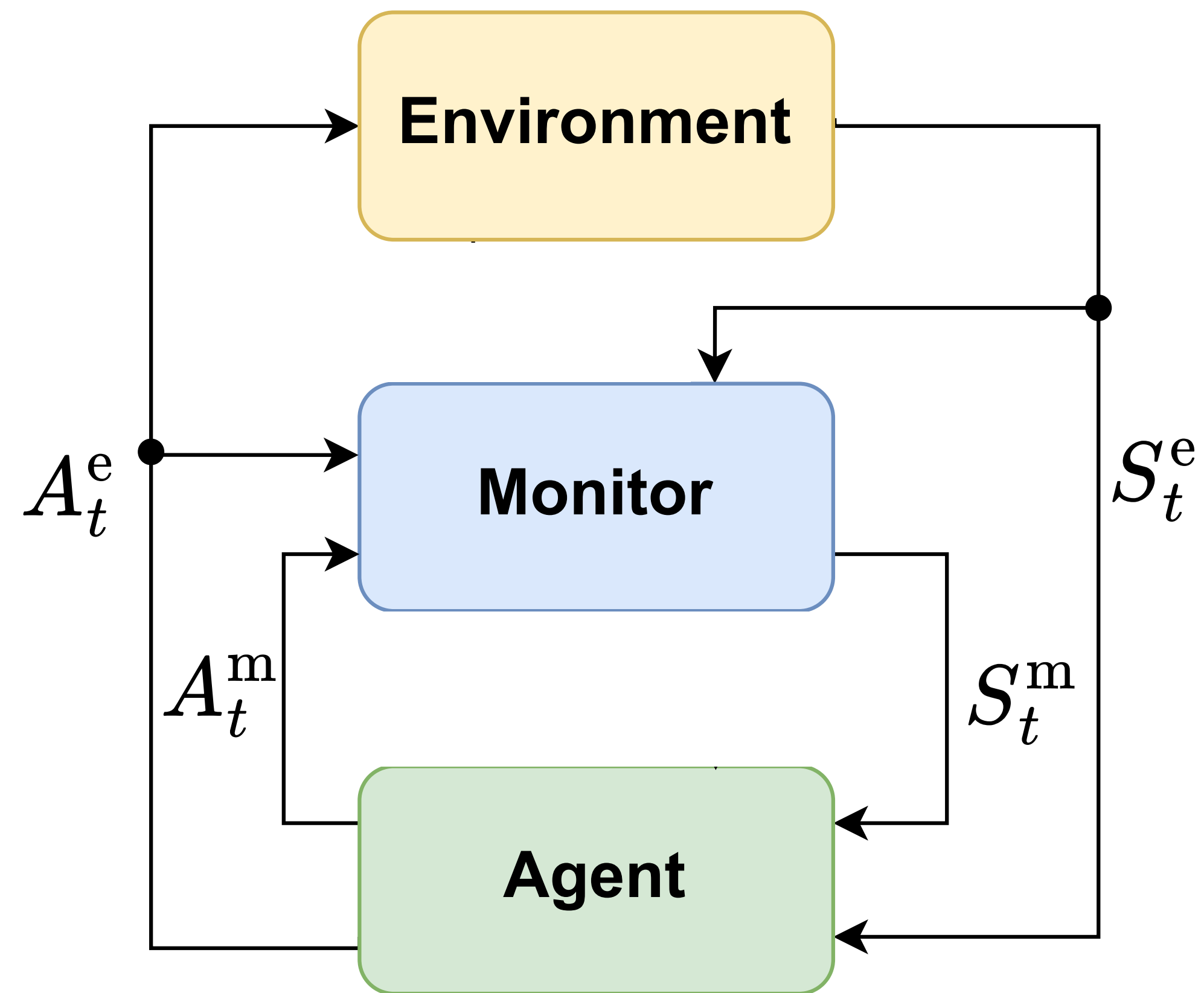
Monitored Markov decision processes (Mon-MDPs)

An extension of MDPs to cover partial observability of rewards



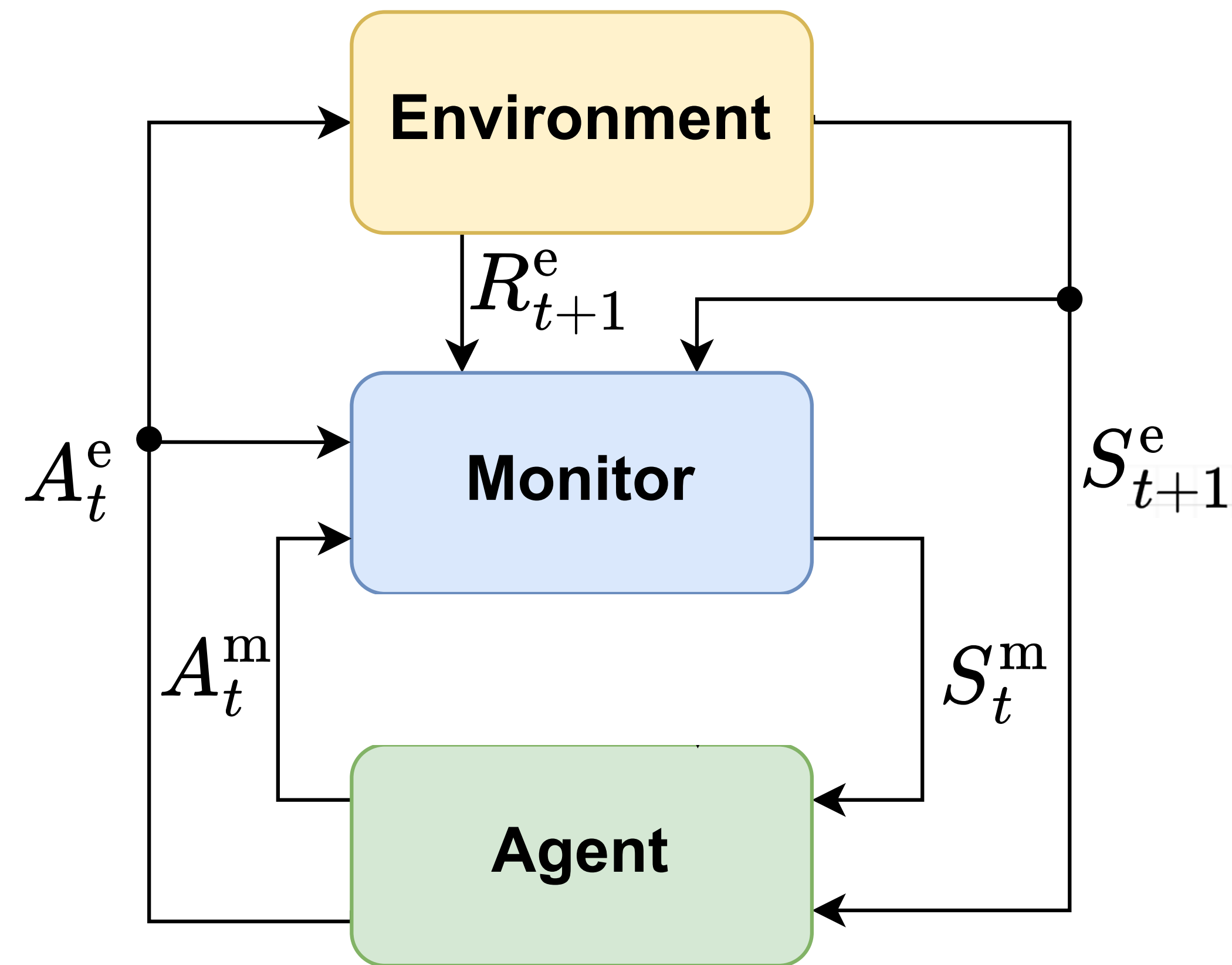
Monitored Markov decision processes (Mon-MDPs)

An extension of MDPs to cover partial observability of rewards



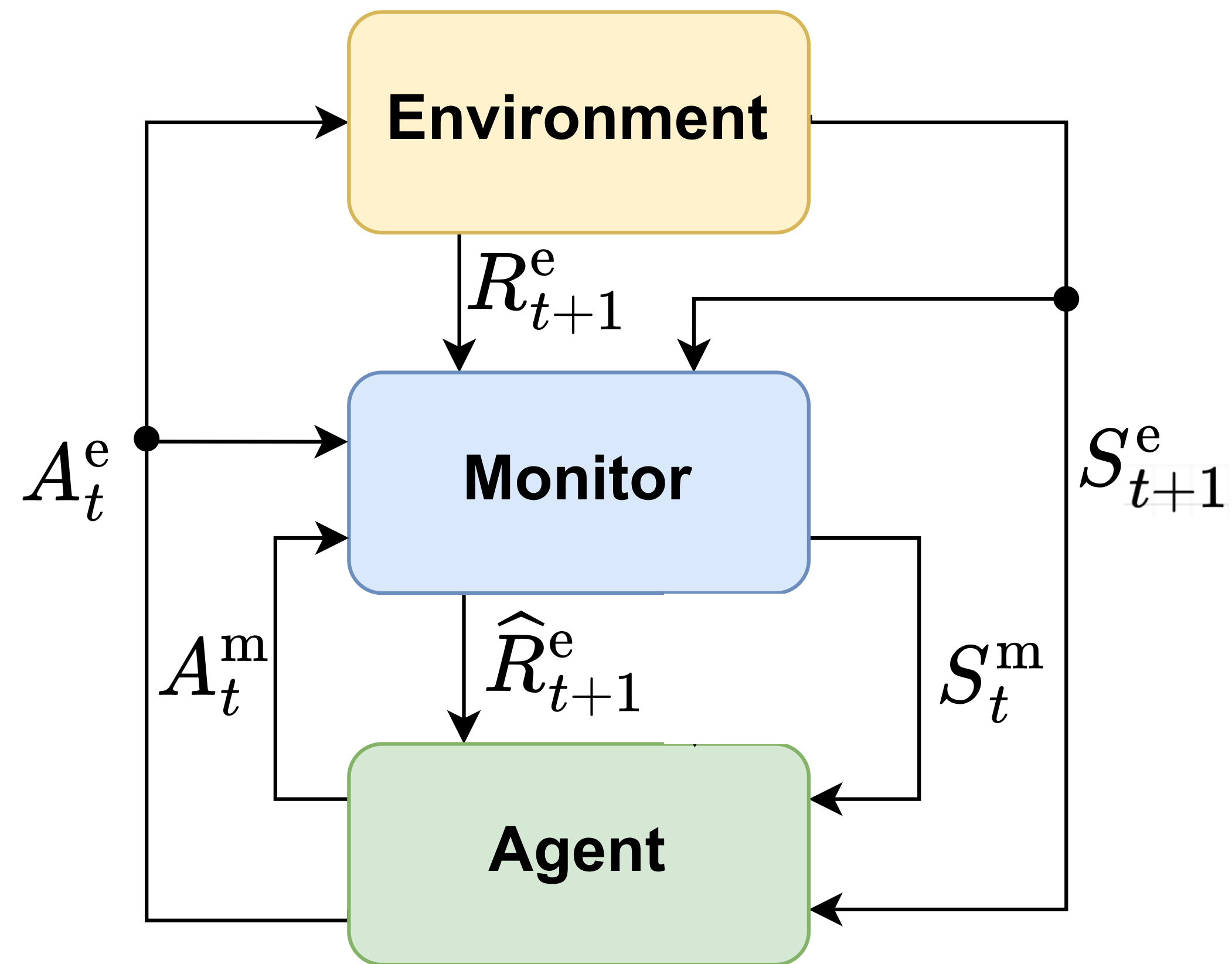
Monitored Markov decision processes (Mon-MDPs)

An extension of MDPs to cover partial observability of rewards



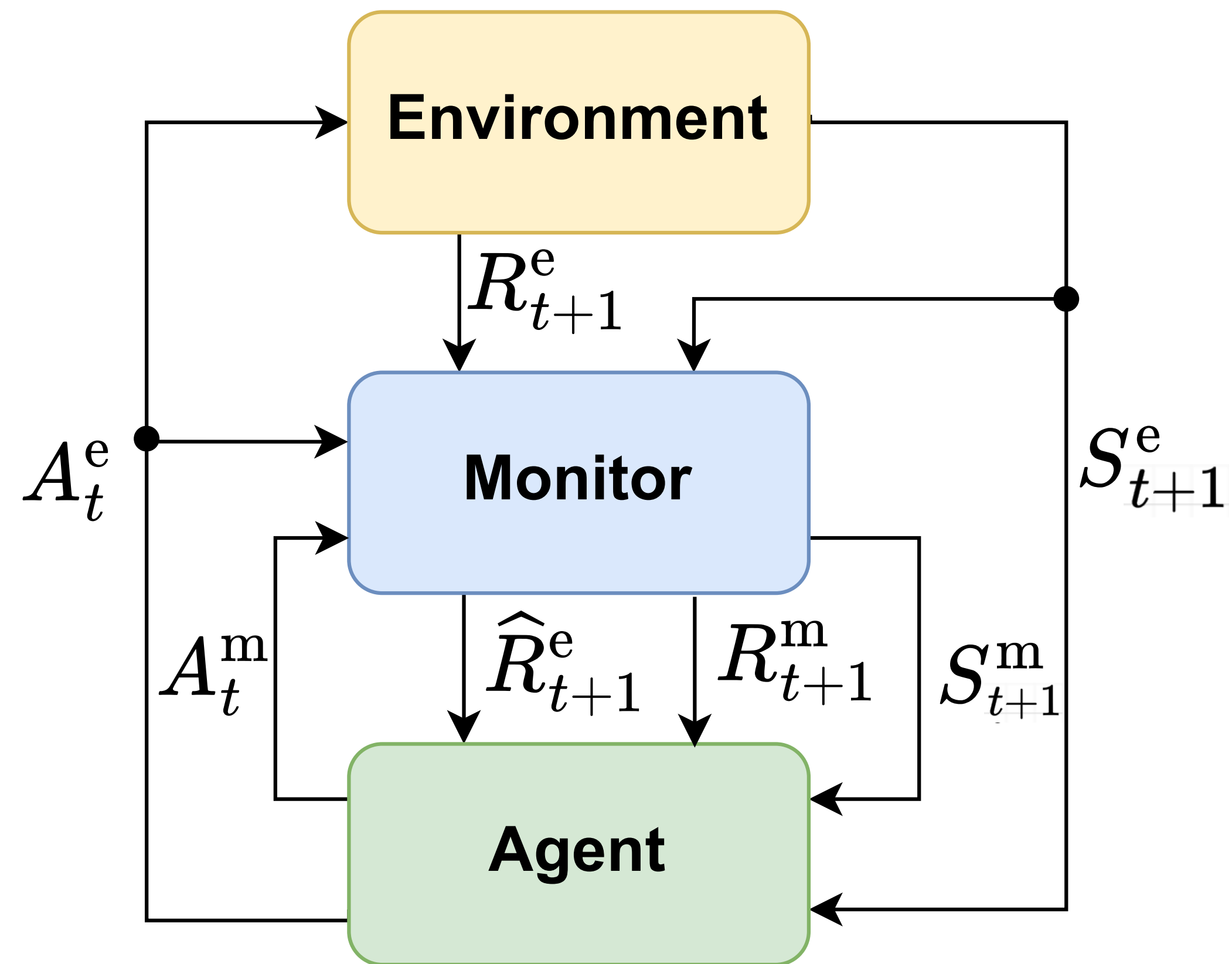
Monitored Markov decision processes (Mon-MDPs)

An extension of MDPs to cover partial observability of rewards



Monitored Markov decision processes (Mon-MDPs)

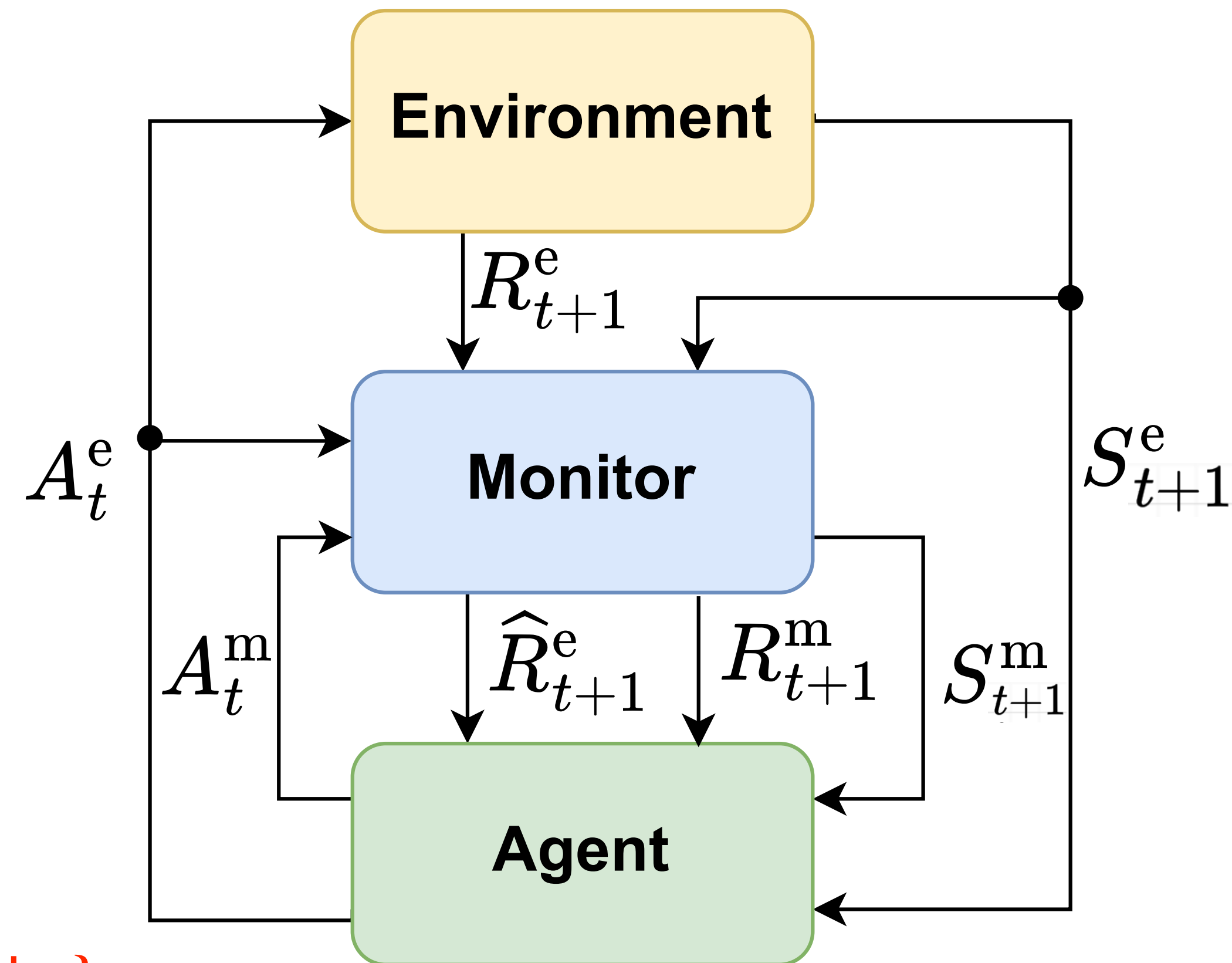
An extension of MDPs to cover partial observability of rewards



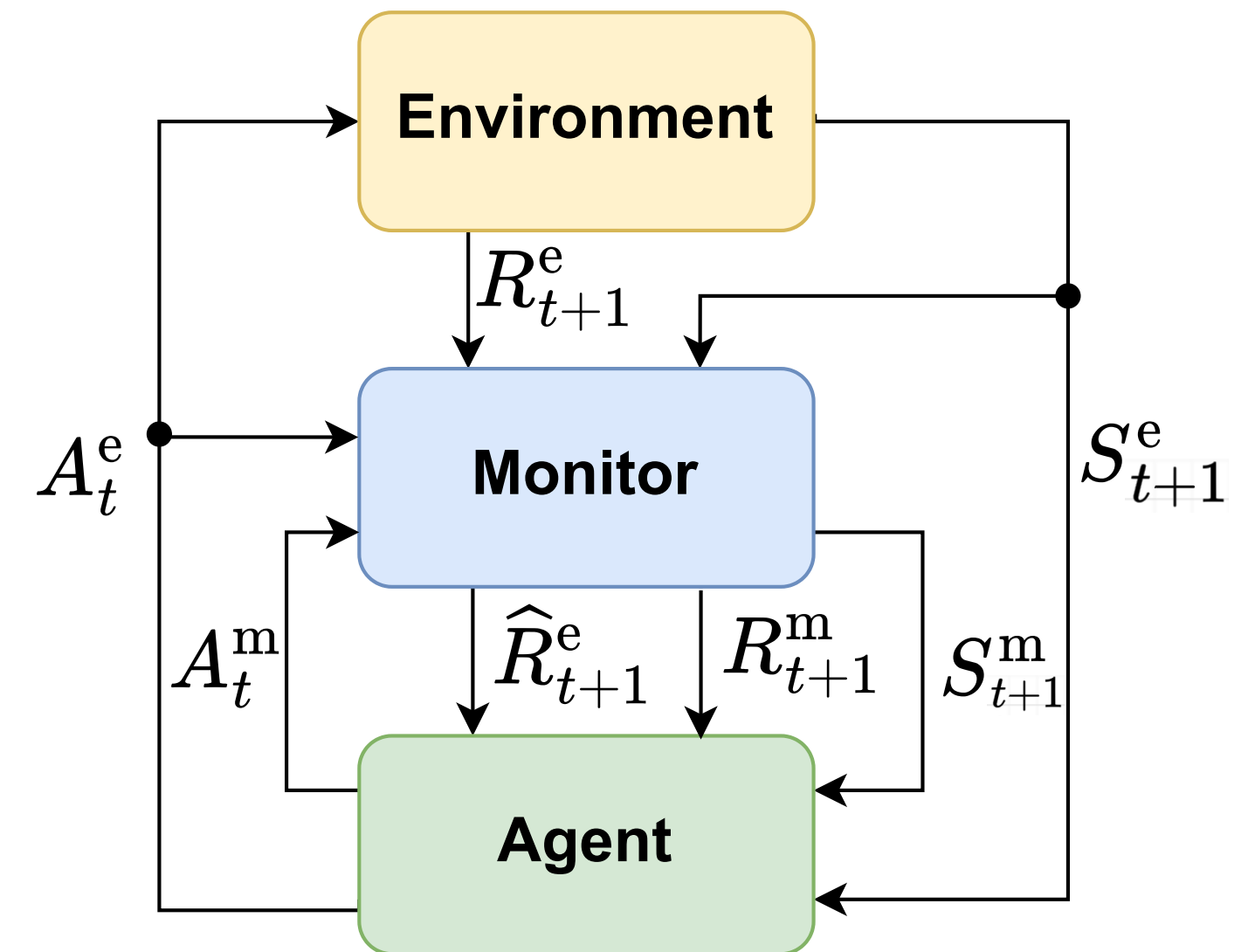
Monitored Markov decision processes (Mon-MDPs)

An extension of MDPs to cover partial observability of rewards

- The goal: $\max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(\textcolor{red}{R}_{t+1}^e + R_{t+1}^m \right) \right]$
- A finite Mon-MDP: $\langle \mathcal{S}, \mathcal{A}, r, p, f^m, \gamma \rangle$
- $\mathcal{S} := \mathcal{S}^e \times \mathcal{S}^m$, $\mathcal{A} := \mathcal{A}^e \times \mathcal{A}^m$
- r is the joint mean reward
- p is the joint transition dynamics
- Monitor function $\hat{R}_{t+1}^e \sim f^m$, and $\hat{R}_{t+1}^e \in \mathbb{R} \cup \{ \textcolor{red}{\perp} \}$
- $0 \leq \gamma < 1$



Assumption

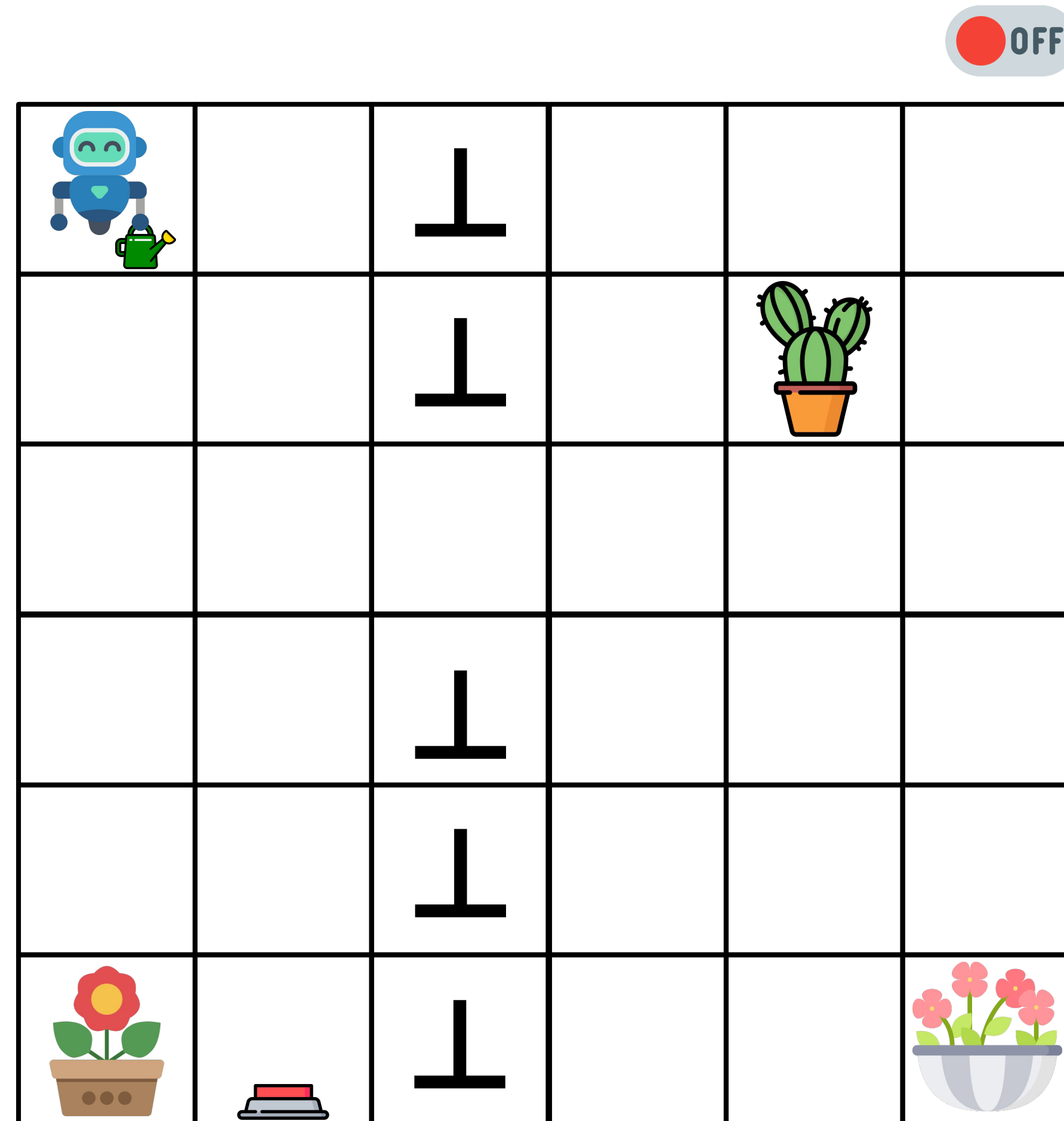


Truthfulness: The monitor doesn't change the underlying reward:

$$\hat{R}_{t+1}^e \in \{R_{t+1}^e, \perp\}$$

Bottleneck - An example of a Mon-MDP

Suppose the button activates a monitoring system



Bottleneck - An example of a Mon-MDP

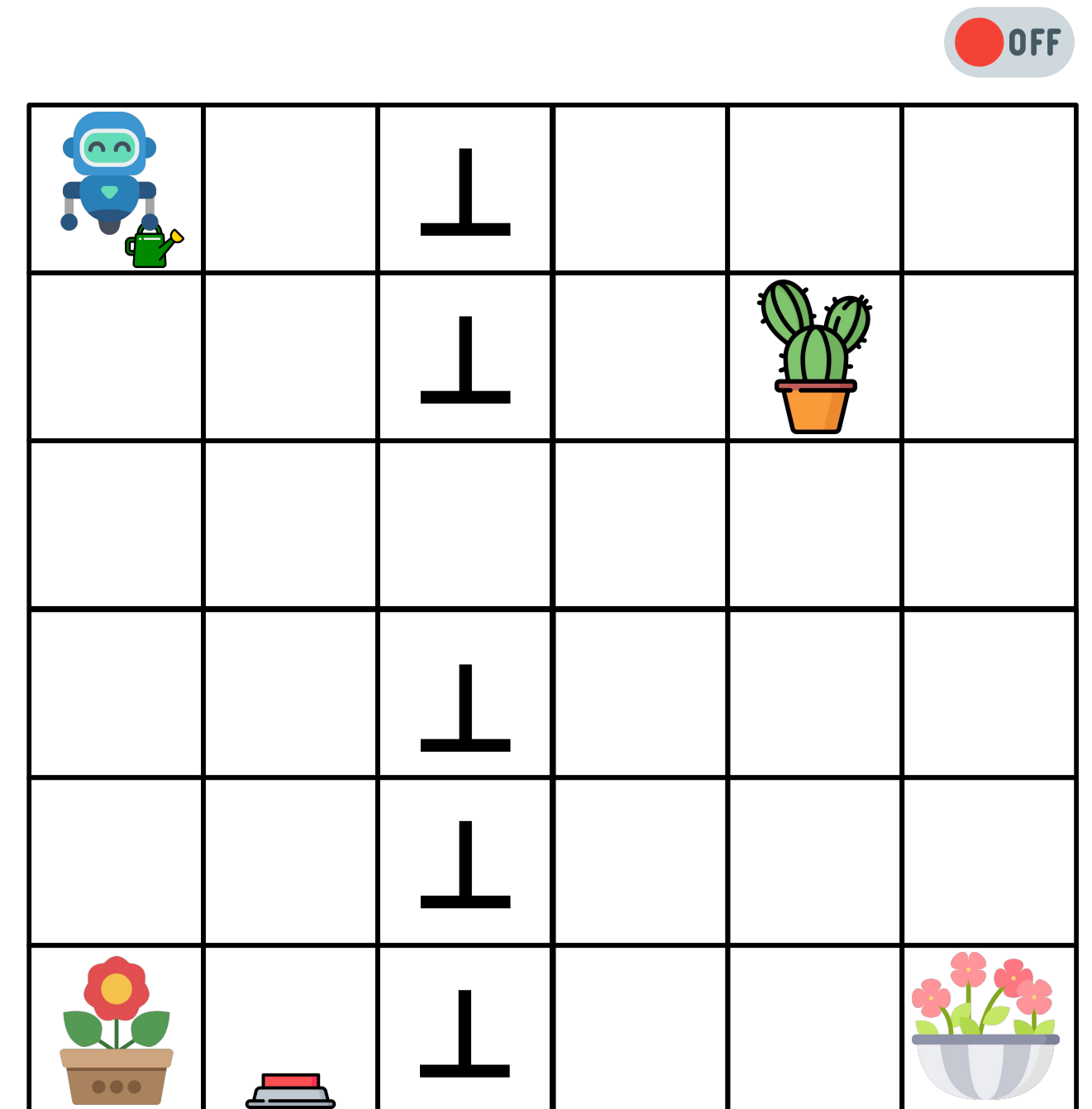
Suppose the button activates a monitoring system

$$\mathcal{S}^m := \{\text{OFF}, \text{ON}\}, \quad \mathcal{A}^m := \{\text{NO-OP}\}$$

Let X_t be random uniform and $0 \leq \rho \leq 1$

$$\widehat{R}_{t+1}^e := \begin{cases} R_{t+1}^e, & \text{if } X_t \leq \rho \text{ and } S_t^m = \text{ON}; \\ \perp, & \text{Otherwise} \end{cases}$$

$$S_{t+1}^m := \begin{cases} \text{ON}, & \text{if } S_t^m = \text{OFF and } S_t^e = \text{"B-CELL"} \text{ and } A_t^e = \downarrow ; \\ \text{OFF}, & \text{if } S_t^m = \text{ON and } S_t^e = \text{"B-CELL"} \text{ and } A_t^e = \downarrow ; \\ S_t^m, & \text{Otherwise} \end{cases}$$



Bottleneck - An example of a Mon-MDP

Suppose the button activates a monitoring system

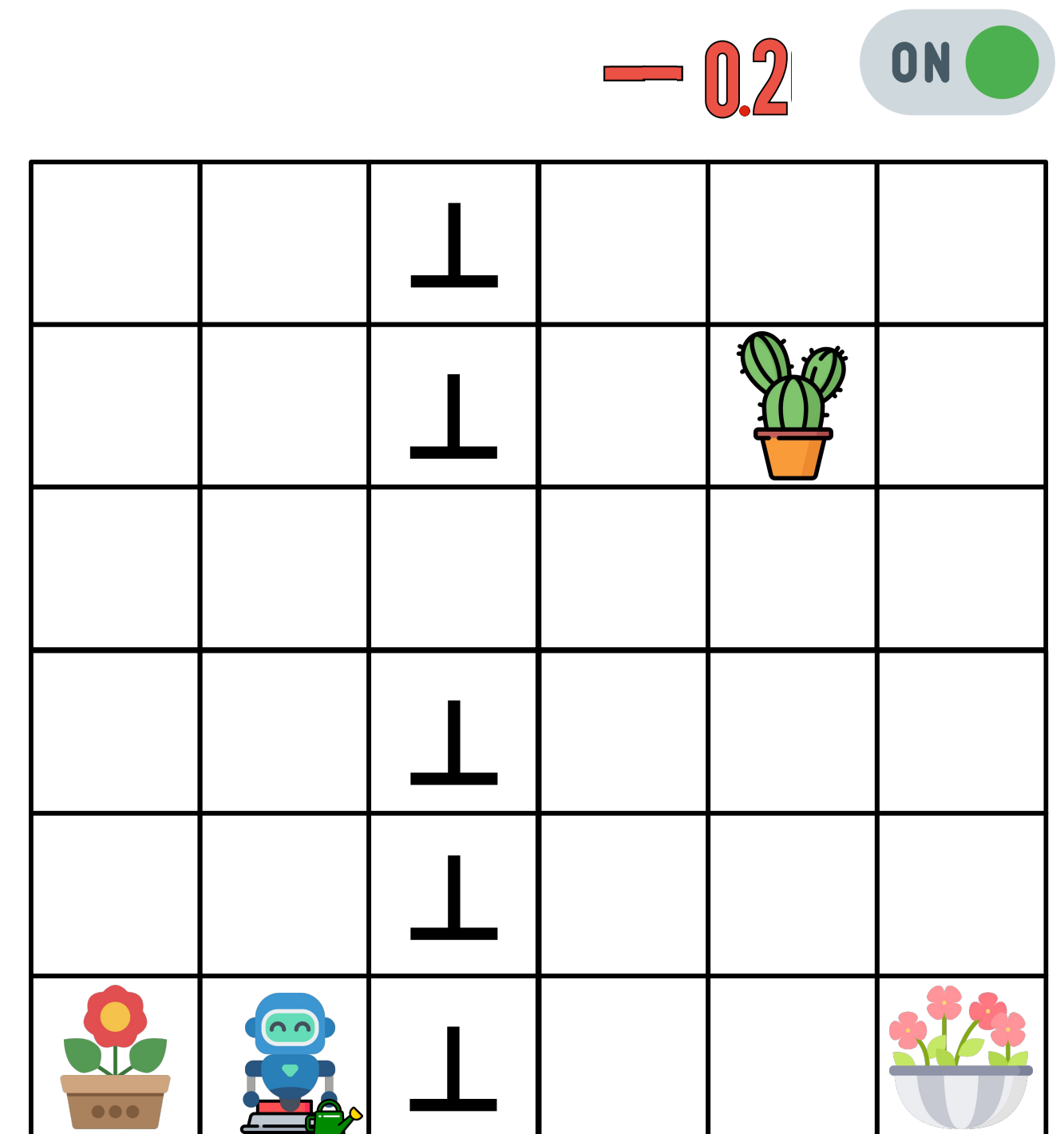
$$\mathcal{S}^m := \{\text{OFF}, \text{ON}\}, \quad \mathcal{A}^m := \{\text{NO-OP}\}$$

Let X_t be random uniform and $0 \leq \rho \leq 1$

$$\widehat{R}_{t+1}^e := \begin{cases} R_{t+1}^e, & \text{if } X_t \leq \rho \text{ and } S_t^m = \text{ON}; \\ \perp, & \text{Otherwise} \end{cases}$$

$$S_{t+1}^m := \begin{cases} \text{ON}, & \text{if } S_t^m = \text{OFF} \text{ and } S_t^e = \text{"B-CELL"} \text{ and } A_t^e = \downarrow; \\ \text{OFF}, & \text{if } S_t^m = \text{ON} \text{ and } S_t^e = \text{"B-CELL"} \text{ and } A_t^e = \downarrow; \\ S_t^m, & \text{Otherwise} \end{cases}$$


$$R_{t+1}^m := \begin{cases} -0.2, & \text{if } S_t^m = \text{ON}; \\ 0, & \text{Otherwise} \end{cases}$$



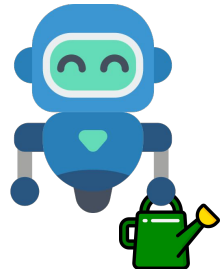

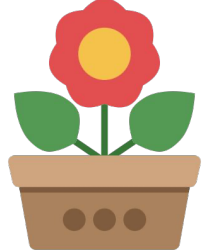


Solution

Our research questions

Review

- 1. How to detect \perp cells from all the others?
- 2. How to deal with \perp cells?
- 3. Can the agent be efficient in watering  while not impacting (1) and (2)?



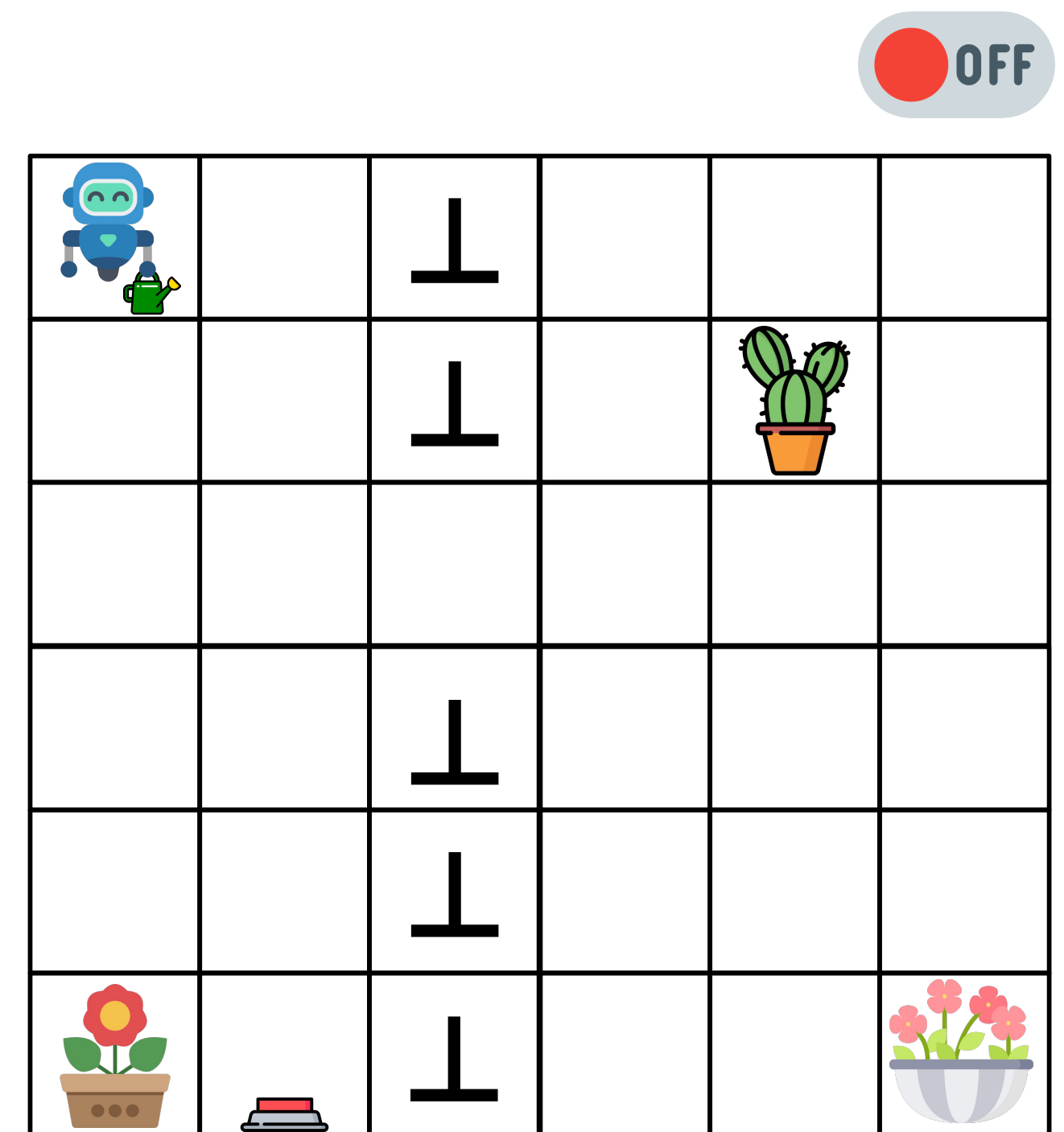
		\perp			
		\perp			
		\perp			
		\perp			
		\perp			

1- How to detect true \perp cells?

1- How to detect true ⊥ cells?

Explore to observe rewards

$$\widetilde{R}_{t+1} = \begin{cases} 1 & \text{if the action led to observing the reward in a state that the reward hasn't been observed before} \\ 0 & \text{otherwise} \end{cases}$$

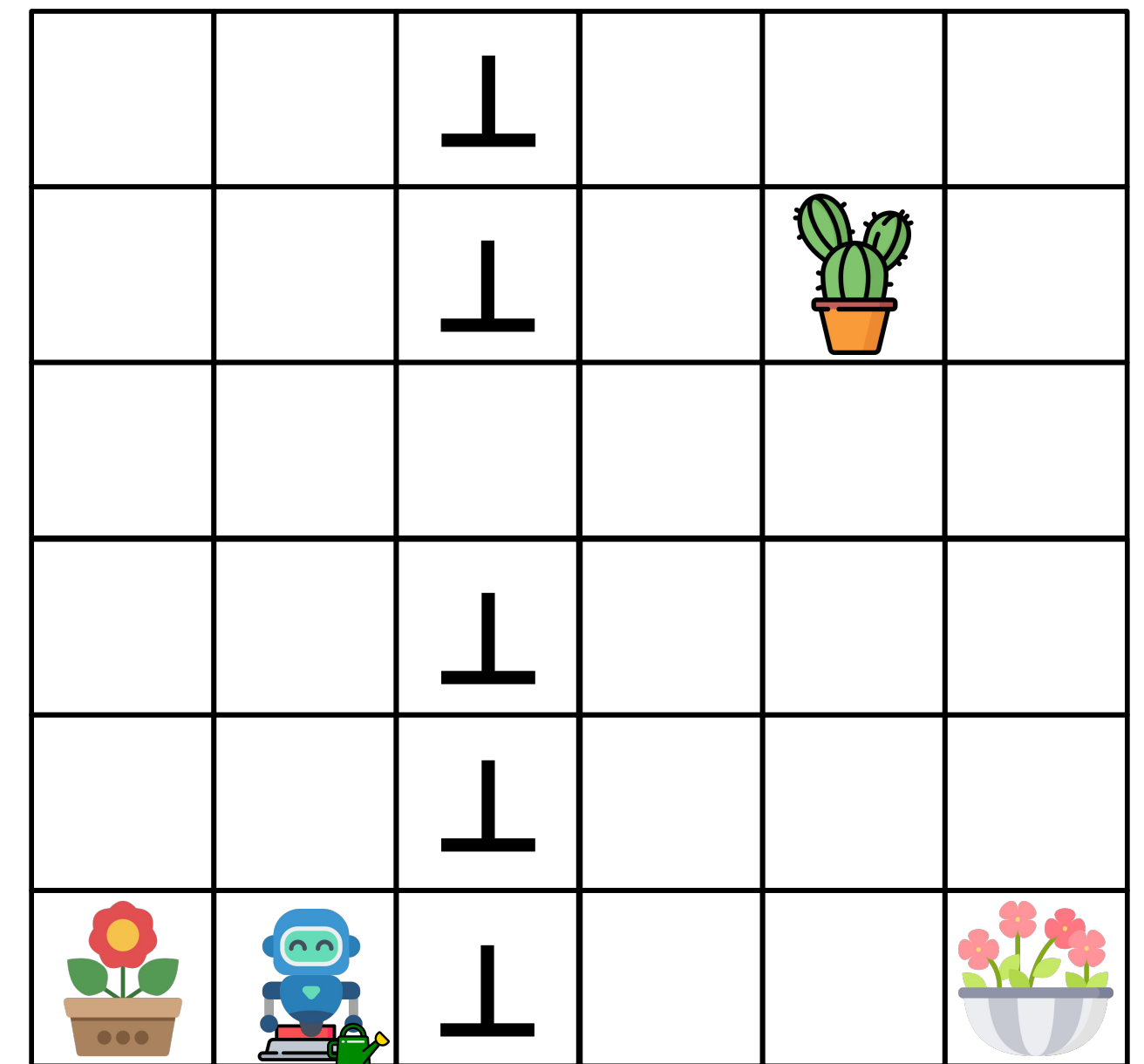


1- How to detect true ⊥ cells?

Explore to observe rewards

$$\widetilde{R}_{t+1} = \begin{cases} 1 & \text{if the action led to observing the reward in a state that the reward hasn't been observed before} \\ 0 & \text{otherwise} \end{cases}$$

ON

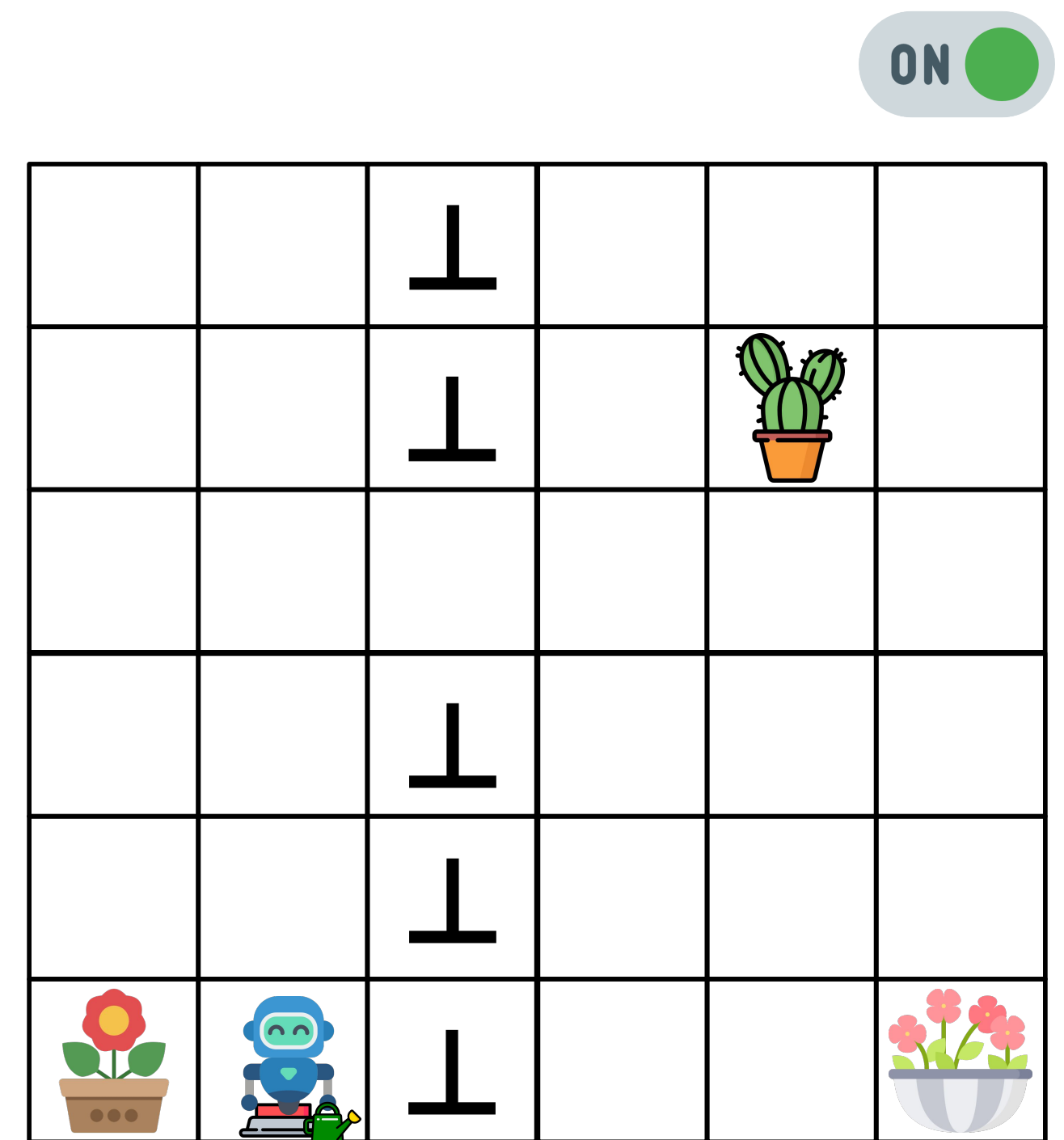


1- How to detect true \perp cells?

Explore to observe rewards

$$\widetilde{R}_{t+1} = \begin{cases} 1 & \text{if the action led to observing the reward in a state that the reward hasn't been observed before} \\ 0 & \text{otherwise} \end{cases}$$

\widetilde{R}_{t+1} is Bernoulli.



Measuring uncertainty

Review

- Suppose you have n samples. Then

$$\text{distance}(\text{empirical mean, true mean}) \leq \frac{\beta}{\sqrt{n}}$$

- If you particularly have n Bernoulli samples. Then

$$\text{distance}(\text{empirical mean, true mean}) \leq \frac{\beta}{n}$$

for sufficiently large value of β

$\frac{\beta}{\sqrt{n}}$, and $\frac{\beta}{n}$ measure the
uncertainty.

1- How to detect true \perp cells?

Explore to observe rewards

$$\widetilde{R}_{t+1} = \begin{cases} 1 & \text{if the action led to observing the reward in a state that the reward hasn't been observed before} \\ 0 & \text{otherwise} \end{cases}$$

\widetilde{R}_{t+1} is Bernoulli.

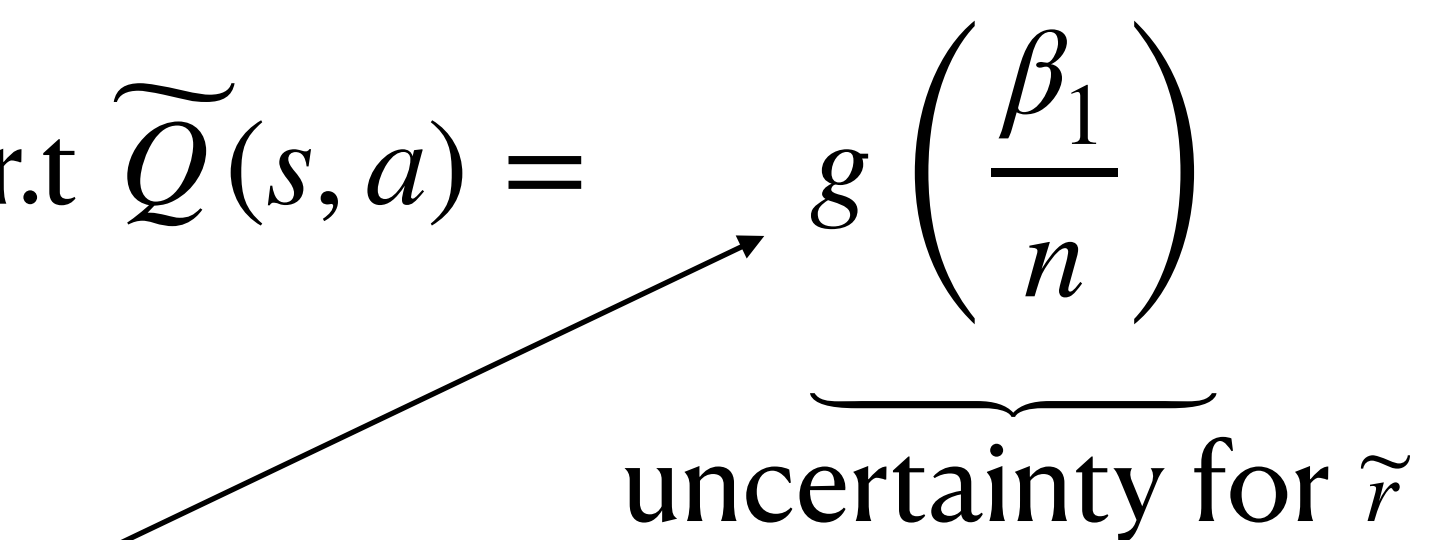
1- How to detect true \perp cells?

Explore to observe rewards

$$\widetilde{R}_{t+1} = \begin{cases} 1 & \text{if the action led to observing the reward in a state that the reward hasn't been observed before} \\ 0 & \text{otherwise} \end{cases}$$

\widetilde{R}_{t+1} is Bernoulli.

Be greedy w.r.t $\widetilde{Q}(s, a) =$


$$g\left(\frac{\beta_1}{n}\right)$$

uncertainty for \tilde{r}

A function

1- How to detect true \perp cells?

Explore to observe rewards

$$\widetilde{R}_{t+1} = \begin{cases} 1 & \text{if the action led to observing the reward in a state that the reward hasn't been observed before} \\ 0 & \text{otherwise} \end{cases}$$

\widetilde{R}_{t+1} is Bernoulli.

Be greedy w.r.t $\widetilde{Q}(s, a) =$


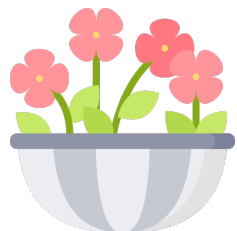
$$\underbrace{g\left(\frac{\beta_1}{n}\right)}_{\text{uncertainty for } \widetilde{r}} + \gamma \sum_{s'} \hat{p}(s' | s, a) \widetilde{V}(s') + \underbrace{\frac{\beta_2}{\sqrt{n}}}_{\text{uncertainty for } \hat{p}}$$

number visits to (s, a)

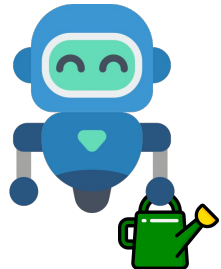




A function

Our research questions

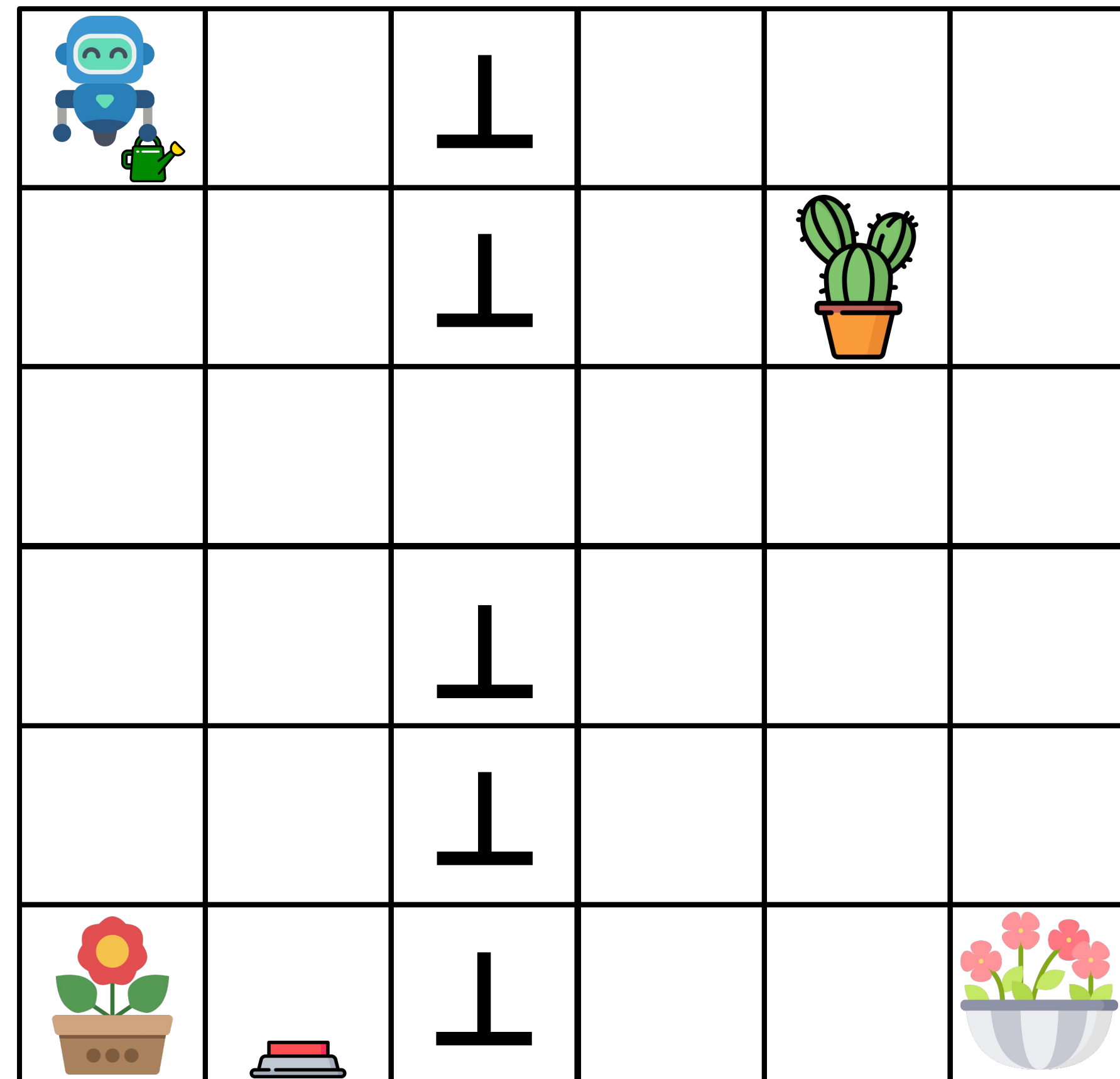
Review

- 1. How to detect \perp cells from all the others? 
- 2. How to deal with \perp cells?
- 3. Can the agent be efficient in watering  while not impacting (1) and (2)?



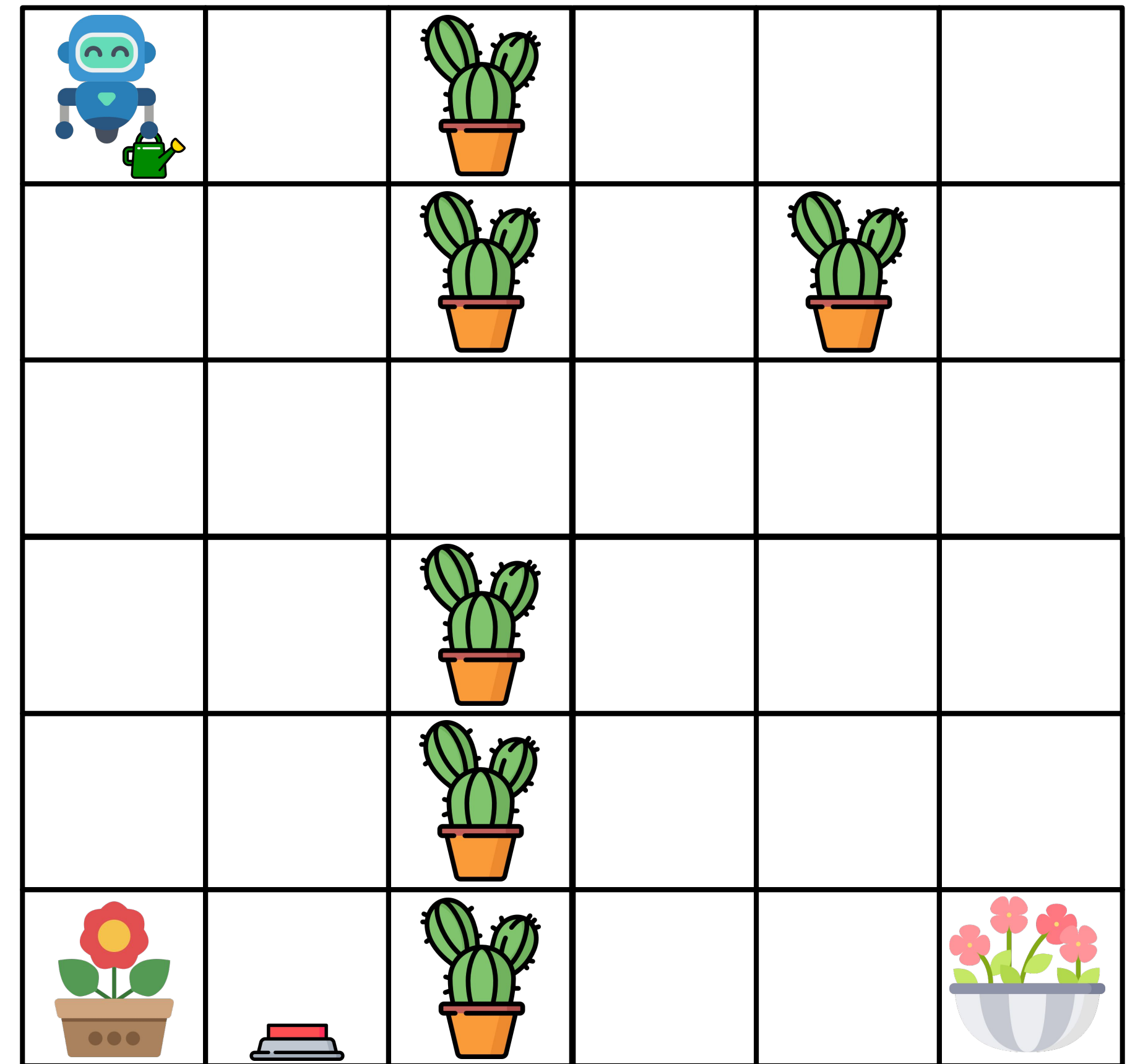
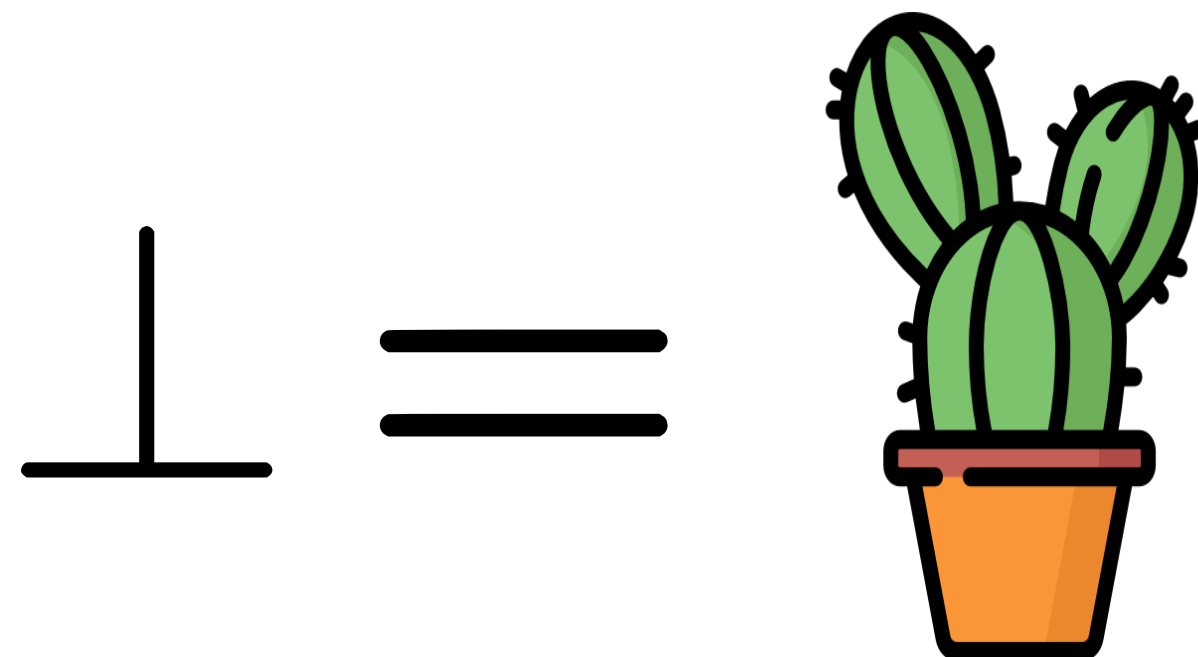
		\perp			
		\perp			
		\perp			
		\perp			
		\perp			

2- How to deal with \perp cells?






2- How to deal with \perp cells?

Be pessimistic about them

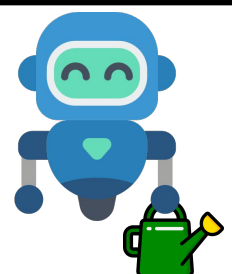
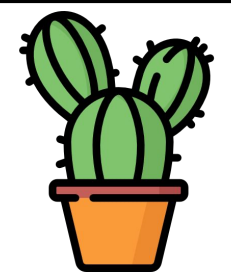
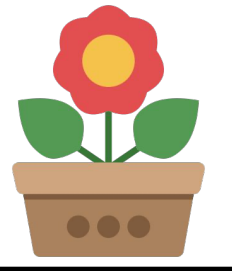




Our research questions

Review

- 1. How to detect \perp cells from all the others? 
- 2. How to deal with \perp cells? 
- 3. Can the agent be efficient in watering  while not impacting (1) and (2)?



		\perp			
		\perp			
		\perp			
		\perp			
		\perp			

3- Can the agent be efficient?

3- Can the agent be efficient?

Use MBIE-EB

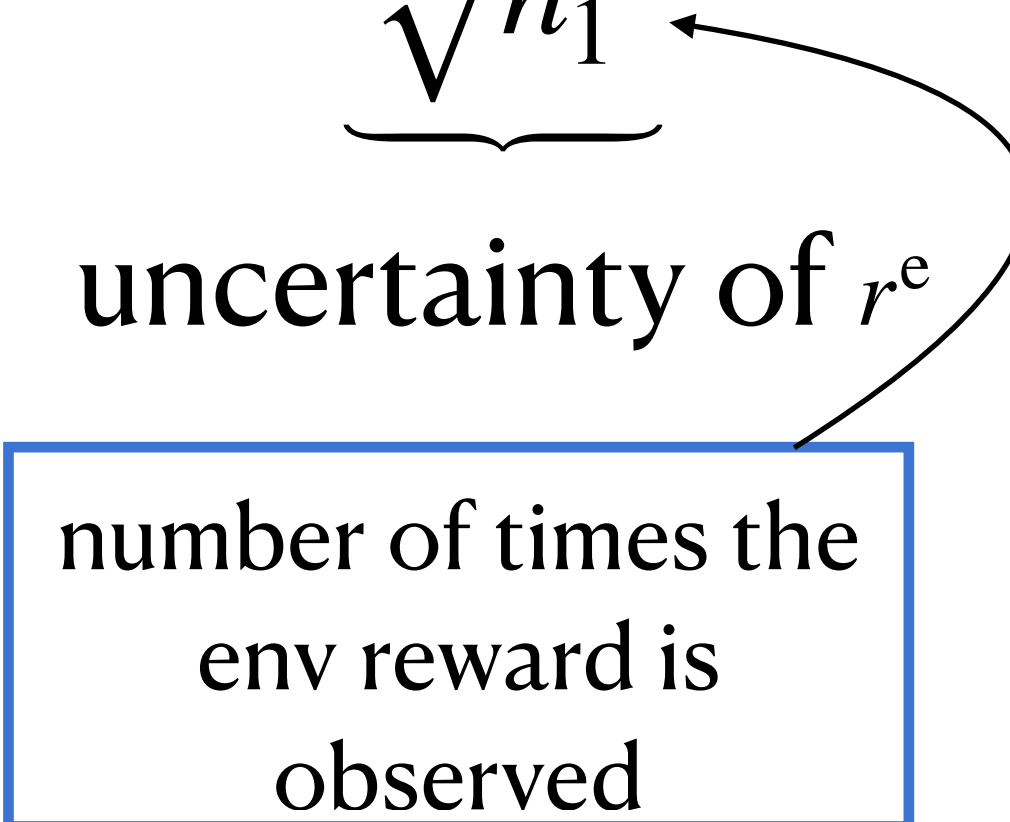
$\frac{\beta}{\sqrt{n}}$, and $\frac{\beta}{n}$ measure the
uncertainty.

3- Can the agent be efficient?

Use MBIE-EB

$$\hat{Q}(s, a) = \hat{r}^e(s^e, a^e) + \underbrace{\frac{\beta_1}{\sqrt{n_1}}}_{\text{uncertainty of } r^e}$$

number of times the
env reward is
observed

A curved arrow points from the text 'number of times the env reward is observed' (which is enclosed in a blue rectangular box) to the denominator $\sqrt{n_1}$ of the uncertainty term in the equation above.

3- Can the agent be efficient?

Use MBIE-EB

$$\hat{Q}(s, a) = \hat{r}^e(s^e, a^e) + \underbrace{\frac{\beta_1}{\sqrt{n_1}}}_{\text{uncertainty of } r^e} + \hat{r}^m(s^m, a^m) + \underbrace{\frac{\beta_2}{\sqrt{n_2}}}_{\text{uncertainty of } r^m}$$

number of times the
env reward is
observed

number of times the
mon reward is
observed

3- Can the agent be efficient?

Use MBIE-EB

$$\hat{Q}(s, a) = \hat{r}^e(s^e, a^e) + \underbrace{\frac{\beta_1}{\sqrt{n_1}}}_{\text{uncertainty of } r^e} + \hat{r}^m(s^m, a^m) + \underbrace{\frac{\beta_2}{\sqrt{n_2}}}_{\text{uncertainty of } r^m} + \gamma \sum_{s'} \hat{p}(s' | s, a) \hat{V}(s') + \underbrace{\frac{\beta_3}{\sqrt{n_3}}}_{\text{uncertainty of } p}$$

The diagram illustrates the components of the MBIE-EB equation and their relationship to observation counts. Three blue-bordered boxes are positioned below the equation, each with an arrow pointing to a specific uncertainty term:

- The first box, labeled "number of times the env reward is observed", points to the term $\frac{\beta_1}{\sqrt{n_1}}$, which is identified as the "uncertainty of r^e ".
- The second box, labeled "number of times the mon reward is observed", points to the term $\frac{\beta_2}{\sqrt{n_2}}$, which is identified as the "uncertainty of r^m ".
- The third box, labeled "number of visits to (s, a) ", points to the term $\frac{\beta_3}{\sqrt{n_3}}$, which is identified as the "uncertainty of p ".

3- Can the agent be efficient?

Use MBIE-EB

$$\hat{Q}(s, a) = \hat{r}^e(s^e, a^e) + \underbrace{\frac{\beta_1}{\sqrt{n_1}}}_{\text{uncertainty of } r^e} + \hat{r}^m(s^m, a^m) + \underbrace{\frac{\beta_2}{\sqrt{n_2}}}_{\text{uncertainty of } r^m} + \gamma \sum_{s'} \hat{p}(s' | s, a) \hat{V}(s') + \underbrace{\frac{\beta_3}{\sqrt{n_3}}}_{\text{uncertainty of } p}$$

number of times the env reward is observed

number of times the mon reward is observed

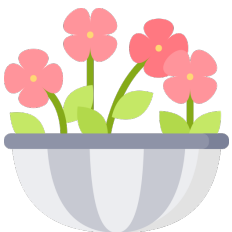
number of visits to (s, a)

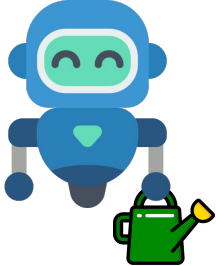
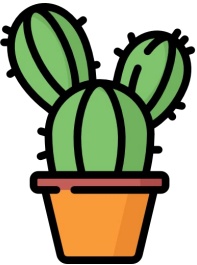
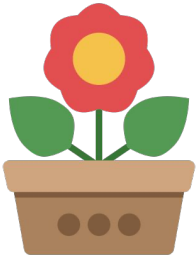


If n_1 was zero (due to unobservability), use  instead.

Our research questions

Review




- 1. How to detect ⊥ cells from all the others? ✓
- 2. How to deal with ⊥ cells? ✓
- 3. Can the agent be efficient in watering  while not impacting (1) and (2)? ✓

		⊥			
		⊥			
		⊥			
		⊥			
		⊥			


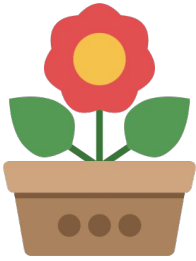

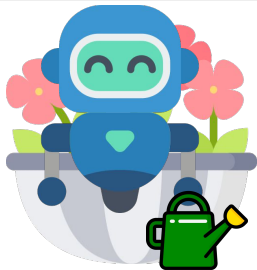
Our research questions

Review



- 1. How to detect ⊥ cells from all the others? ✓
- 2. How to deal with ⊥ cells? ✓
- 3. Can the agent be efficient in watering  while not impacting (1) and (2)? ✓



		⊥			
		⊥			
		⊥			
		⊥			
		⊥			

Monitored MBIE-EB's theoretical performance

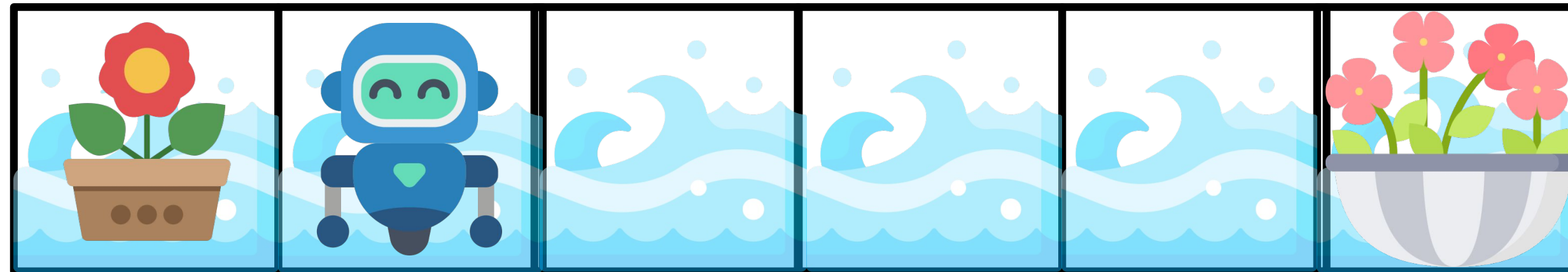
Monitored MBIE-EB's theoretical performance

$$\mathcal{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{\rho(1-\gamma)^6\epsilon^3}\right)$$

Monitored MBIE-EB's empirical performance

Monitored MBIE-EB's empirical performance

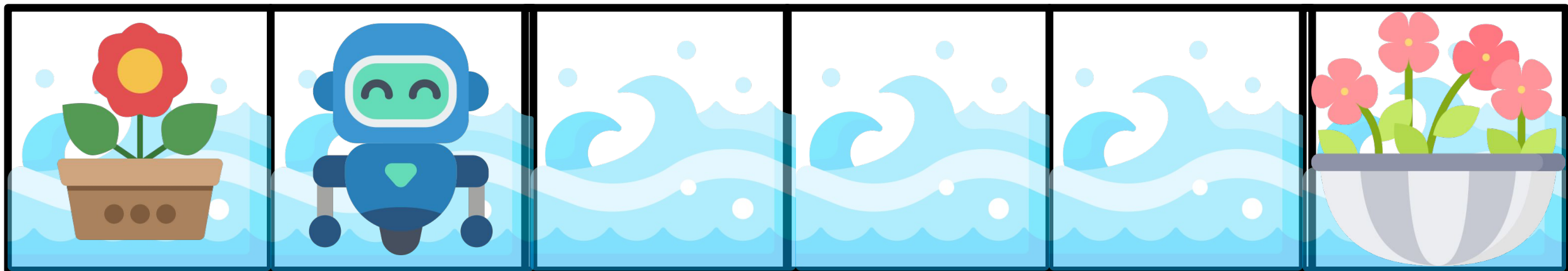
On River Swim



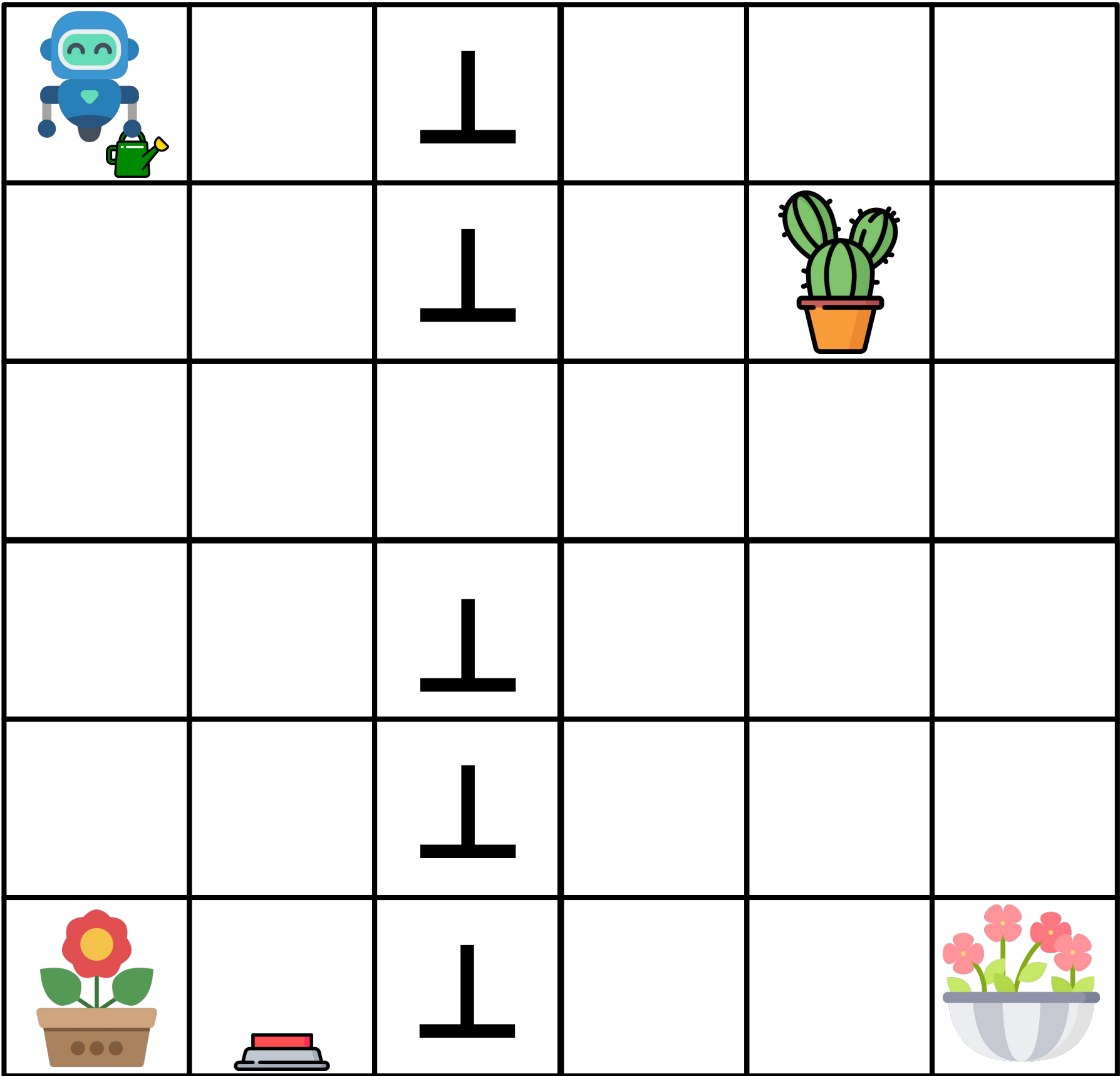
The agent should go to the right but due to stochasticity, it's more likely to move left or stay put. This stochasticity makes the exploration hard.

Monitored MBIE-EB's empirical performance

On River Swim & Bottleneck

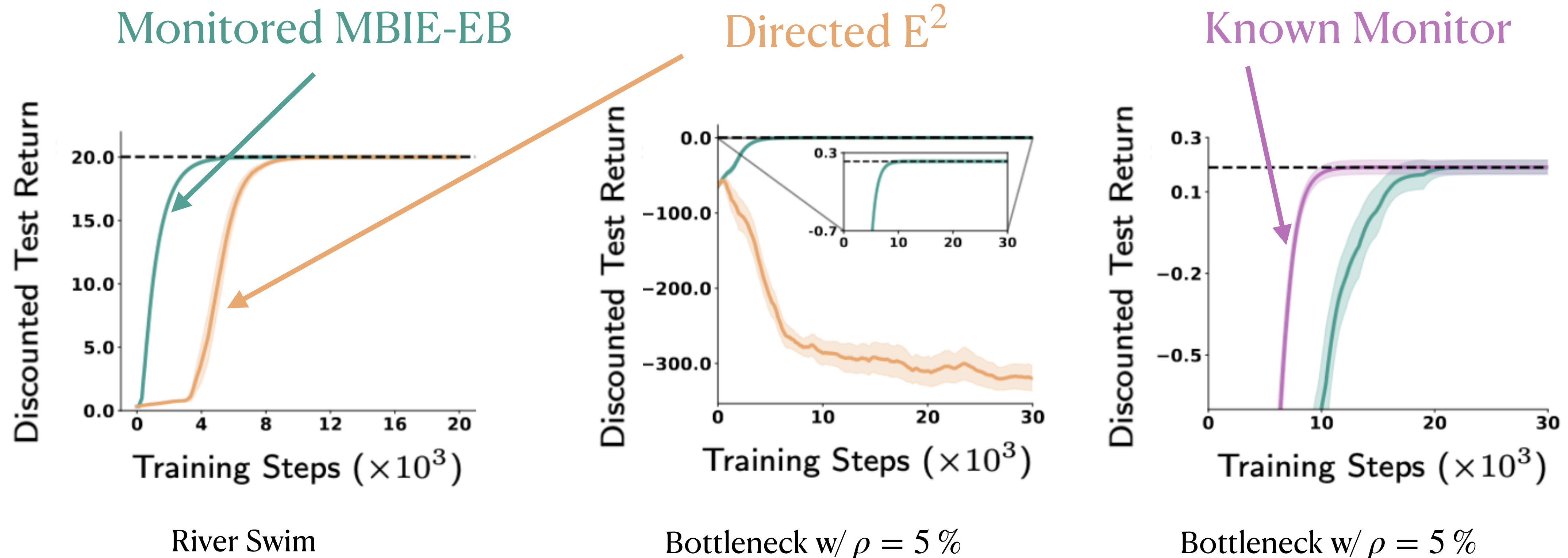


The agent should go to the right but due to stochastic, it's more likely to move left or stay put. This stochasticity makes the exploration hard.



Monitored MBIE-EB's empirical performance

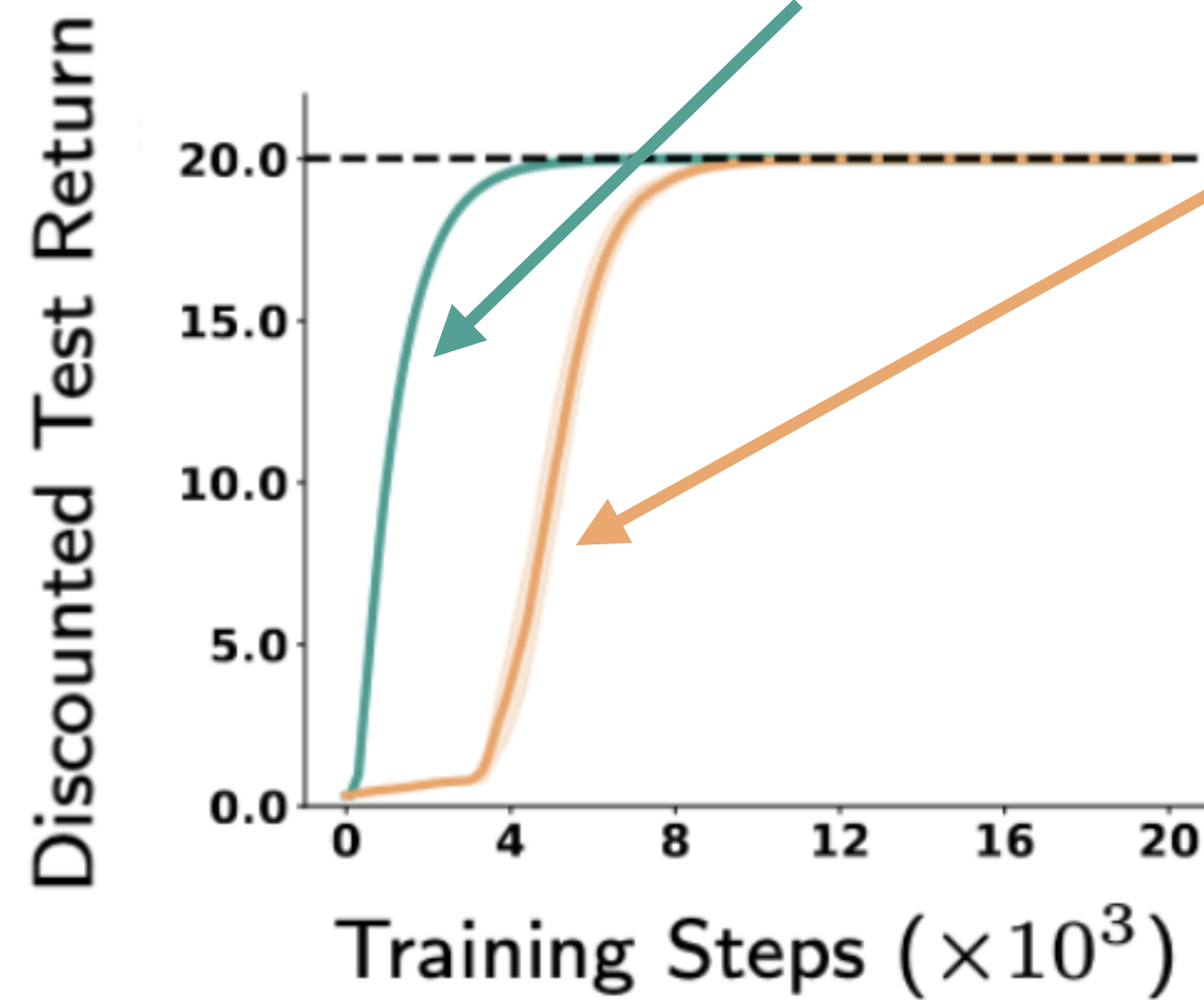
Directed Explore-Exploit (Directed E^2) is the state-of-the-art algorithm in Mon-MDPs



- Dashed horizontal line is the minimax-optimal discounted return.

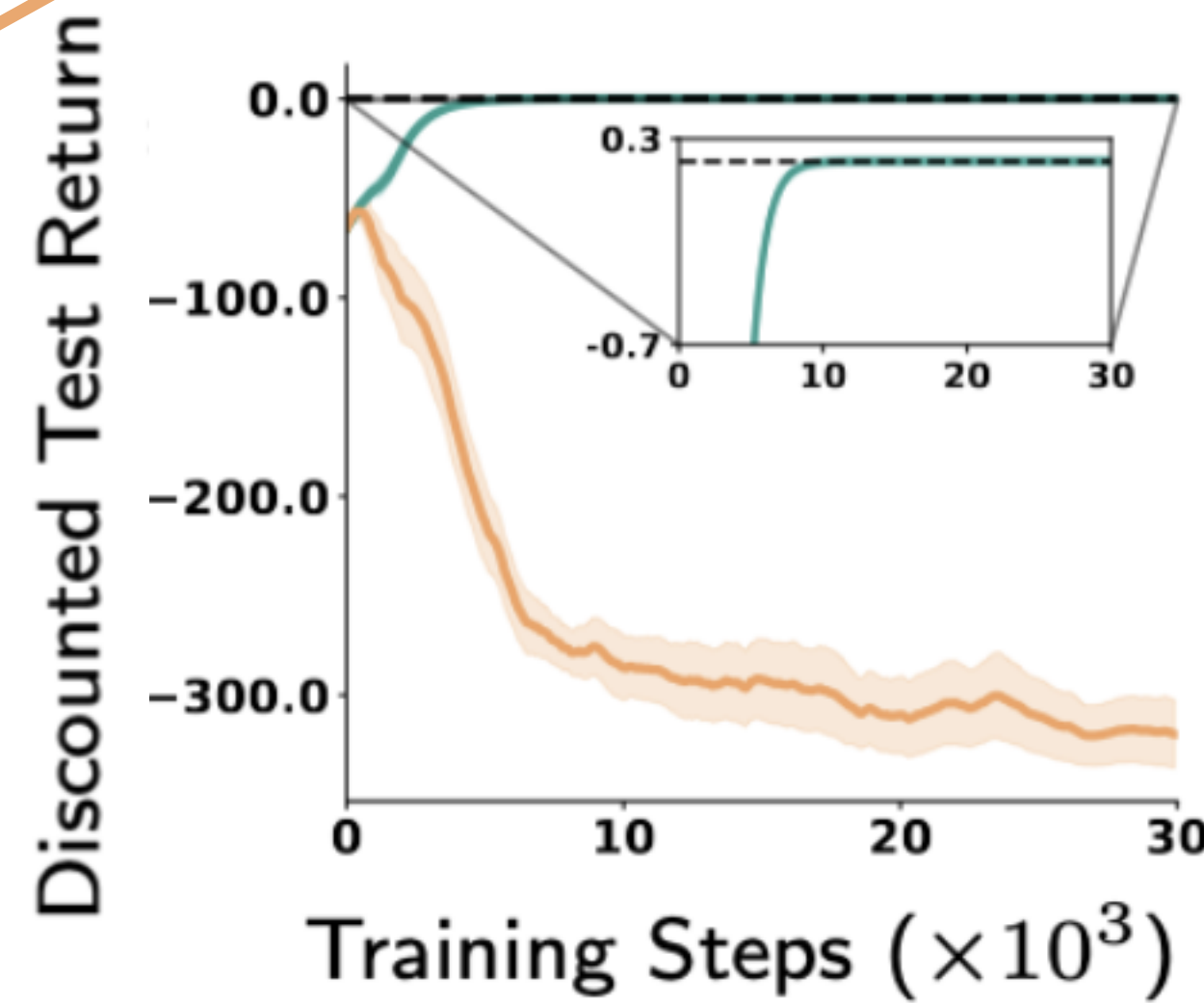
Takeaways

Monitored MBIE-EB



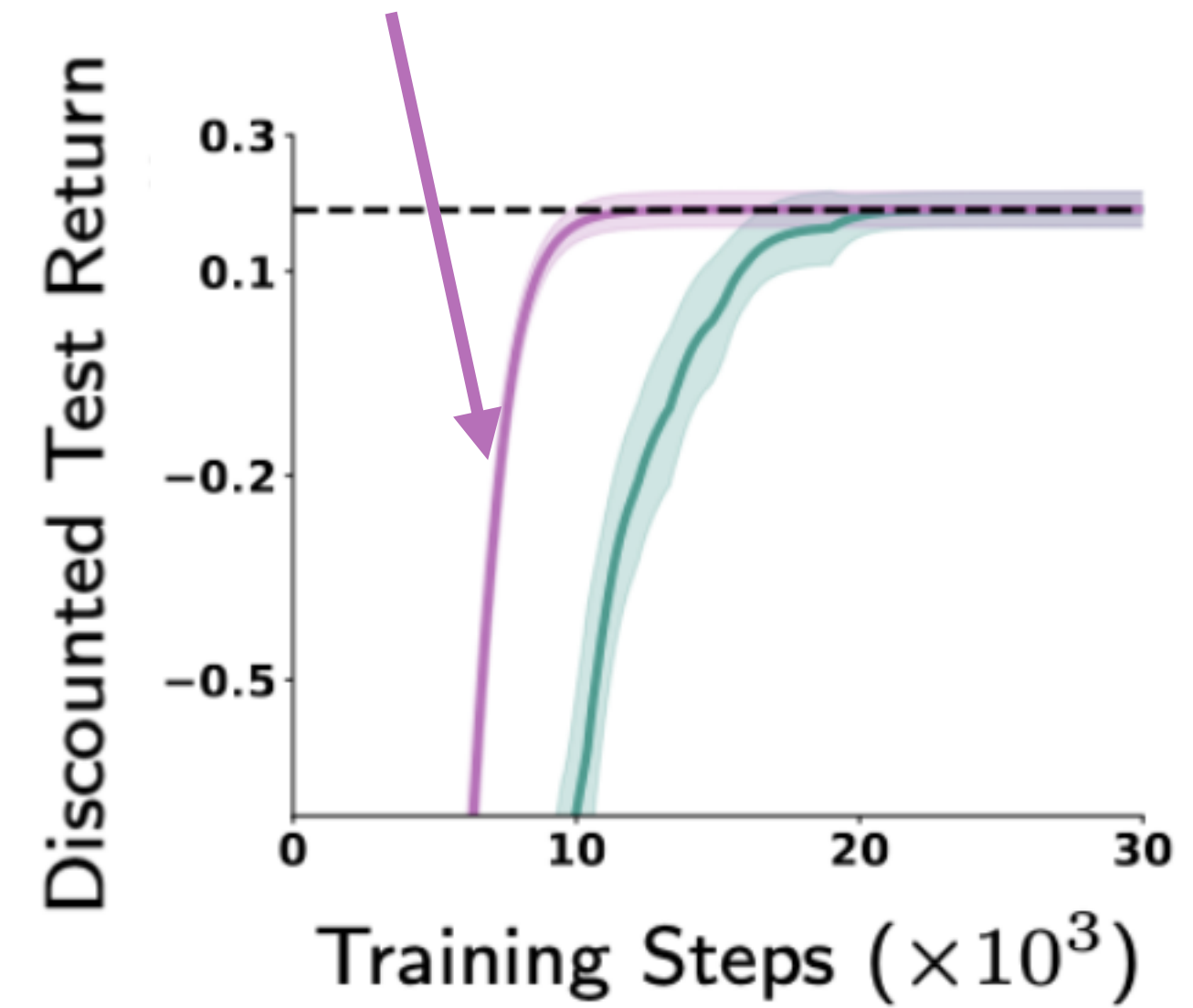
River Swim

Directed E^2



Bottleneck w/ $\rho = 5\%$

Known Monitor



Bottleneck w/ $\rho = 5\%$

1. Due to being model-based and planning, Monitored MBIE-EB performs well on River Swim.
2. Monitored MBIE-EB is robust against stochastic observability and finds the minimax-optimal policy.
3. Monitored MBIE-EB can leverage prior knowledge about the monitor.

List of Contributions

List of my contributions

- Defining the minimax-optimality in Mon-MDPs replacing the notion of MDPs' optimality.
- Presenting Monitored MBIE-EB, the first model-based minimax-optimal algorithm for Mon-MDPs.
- Proving the polynomial sample complexity of Monitored MBIE-EB.
- Showing the dependence of the Monitored MBIE-EB's sample complexity on ρ in Mon-MDPs is essentially unimprovable.
- Demonstrating the superior performance of Monitored MBIE-EB compared to Directed E^2 , the previous state-of-the-art algorithm in Mon-MDPs. We showed more dramatic results when the dynamics of how the agent can or cannot observe the reward is known apriori.

Future work

Beyond finite domains

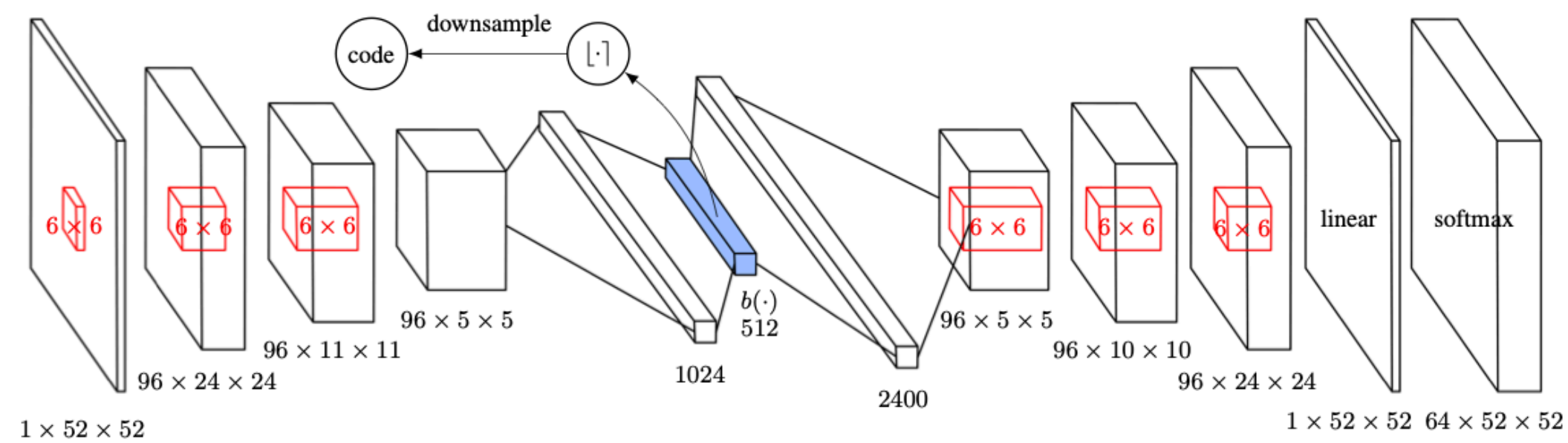
#Exploration: A Study of Count-Based Exploration for Deep Reinforcement Learning

**Haoran Tang^{1*}, Rein Houthoofd^{34*}, Davis Foote², Adam Stooke², Xi Chen^{2†},
Yan Duan^{2†}, John Schulman⁴, Filip De Turck³, Pieter Abbeel^{2†}**

¹ UC Berkeley, Department of Mathematics

² UC Berkeley, Department of Electrical Engineering and Computer Sciences

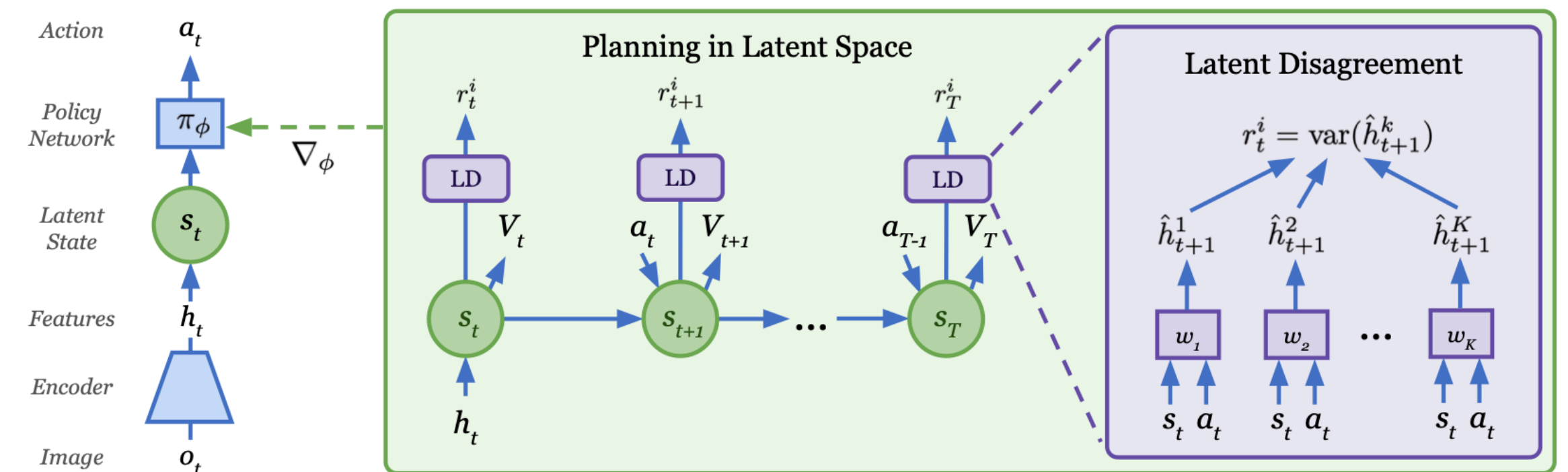
³ Ghent University – imec, Department of Information Technology

⁴ OpenAI

(2016)

Planning to Explore via Self-Supervised World Models

Ramanan Sekar^{1*} Oleh Rybkin^{1*} Kostas Daniilidis¹ Pieter Abbeel² Danijar Hafner^{3,4} Deepak Pathak^{5,6}



(2020)

2- Use a better base algorithm

MBIE-EB's upper bound is loose

Our upper bound

PAC Bounds for Discounted MDPs

Tor Lattimore¹ and Marcus Hutter^{1,2,3}

Research School of Computer Science

¹Australian National University and ²ETH Zürich and ³NICTA
{tor.lattimore,marcus.hutter}@anu.edu.au

$$\tilde{\mathcal{O}} \left(\frac{|\mathcal{S}| |\mathcal{A}|}{\rho(1-\gamma)^6 \epsilon^3} \right)$$

$$\tilde{\Omega} \left(\frac{|\mathcal{S}| |\mathcal{A}|}{(1-\gamma)^3 \epsilon^2} \right)$$

3- Unifying the observation and optimization

A unified algorithm

Near-Optimal Reinforcement Learning in Polynomial Time

MICHAEL KEARNS*

mkearns@cis.upenn.edu

*Department of Computer and Information Science, University of Pennsylvania, Moore School Building,
200 South 33rd Street, Philadelphia, PA 19104-6389, USA*

SATINDER SINGH*

satinder.baveja@syntekcapital.com

Syntek Capital, New York, NY 10019, USA

Explicit Explore or Exploit (E^3), (2002)

R-MAX – A General Polynomial Time Algorithm for Near-Optimal Reinforcement Learning

Ronen I. Brafman

BRAFMAN@CS.BGU.AC.IL

*Computer Science Department
Ben-Gurion University
Beer-Sheva, Israel 84105*

Moshe Tennenholtz*

MOSHE@ROBOTICS.STANFORD.EDU

*Computer Science Department
Stanford University
Stanford, CA 94305*

R-Max, (2002)

Acknowledgement

I can't thank enough these individuals

- My lovely supervisors: Matt and Mike
- My collaborators: Simone and Monta
- Mike's and Matt's lab: Too many to name individually
- My teachers and all theoreticians throughout history whose contributions has given and will give me insights to study and (ideally) solve problems
- My family
- God ❤️

Thank you! :)