

# OUTCOME-DRIVEN REINFORCEMENT LEARNING VIA VARIATIONAL INFERENCE



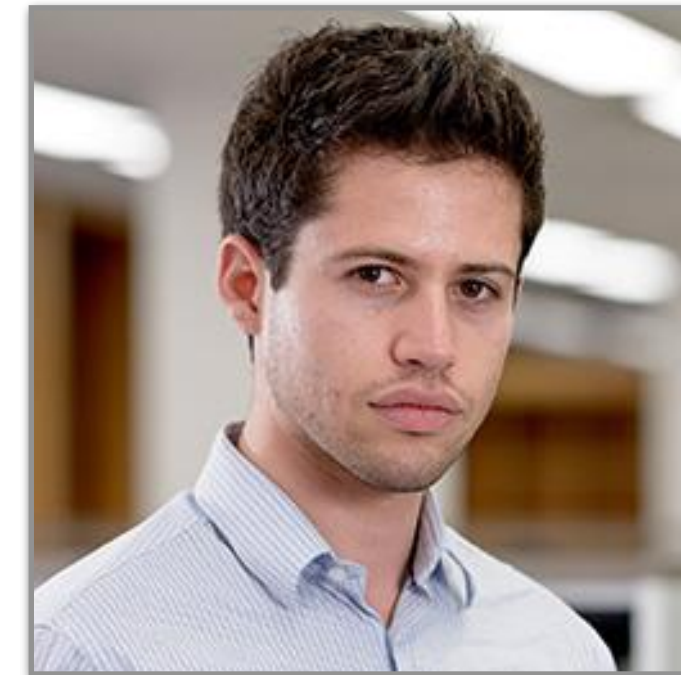
**TIM G. J. RUDNER\***



**VITCHYR H. PONG\***



ROWAN McALLISTER



YARIN GAL



SERGEY LEVINE

Neural Information Processing Systems 2021



Correspondence to

`tim.rudner@cs.ox.ac.uk`

`vitchy@berkeley.edu`



**Can we derive RL from probabilistic inference?**

**Can we derive RL from probabilistic inference  
without access to a reward function?**

**Yes!**

**Yes!**

**A reward function**

**Yes!**

**A reward function **emerges naturally** from inference,**

**Yes!**

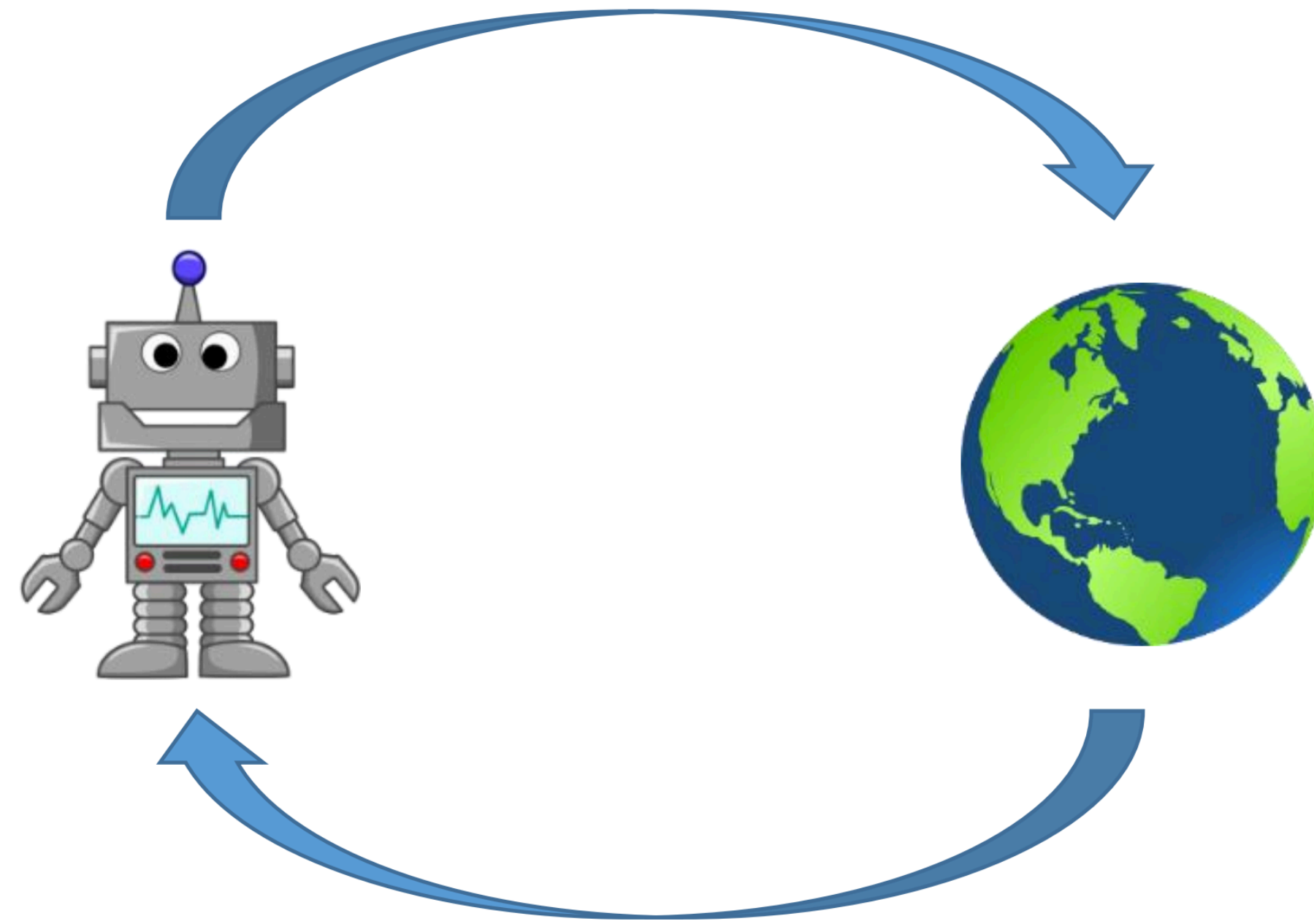
**A reward function emerges naturally from inference,  
can be learned from data, and**

**Yes!**

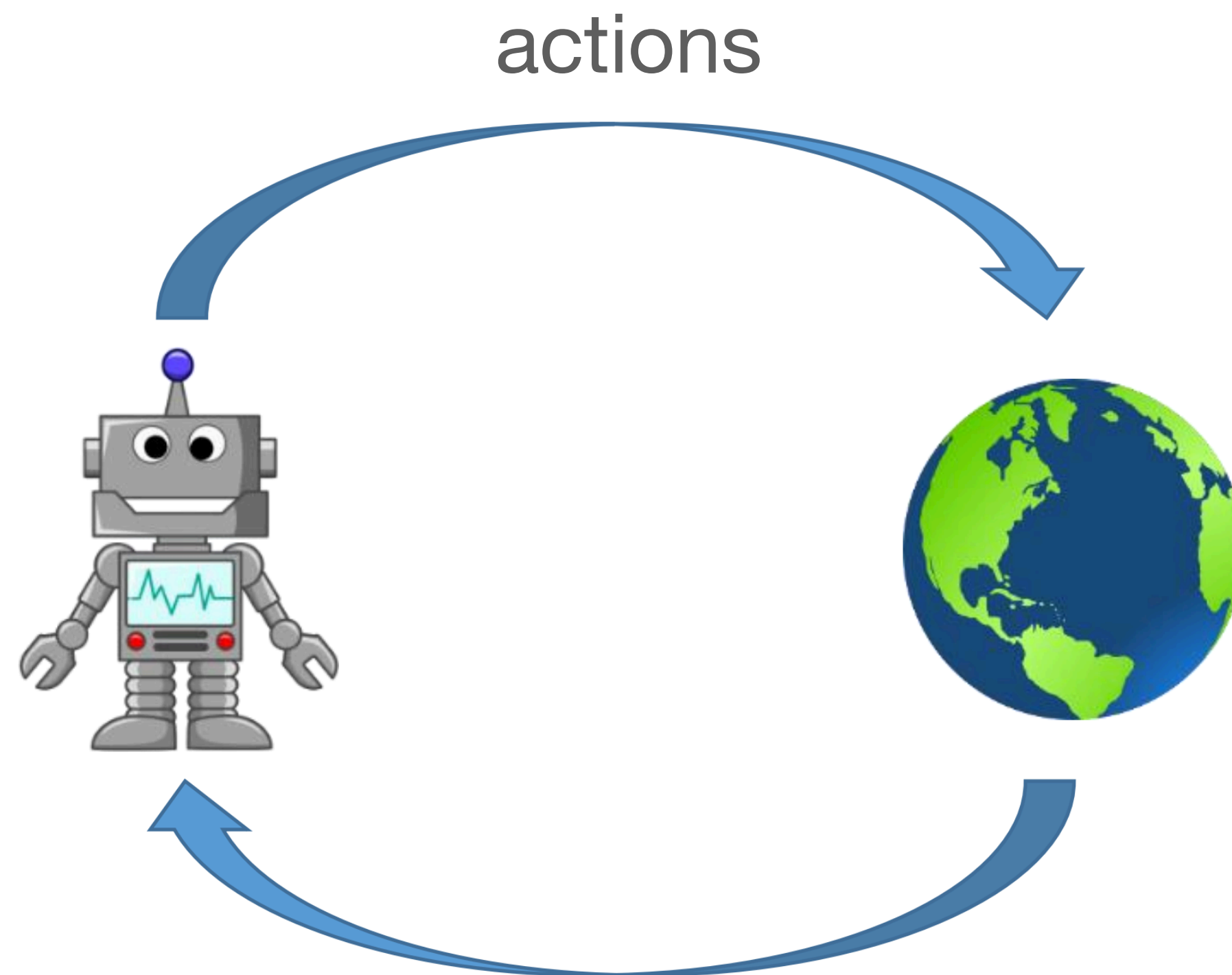
**A reward function emerges naturally from inference,  
can be learned from data, and  
is well-shaped.**



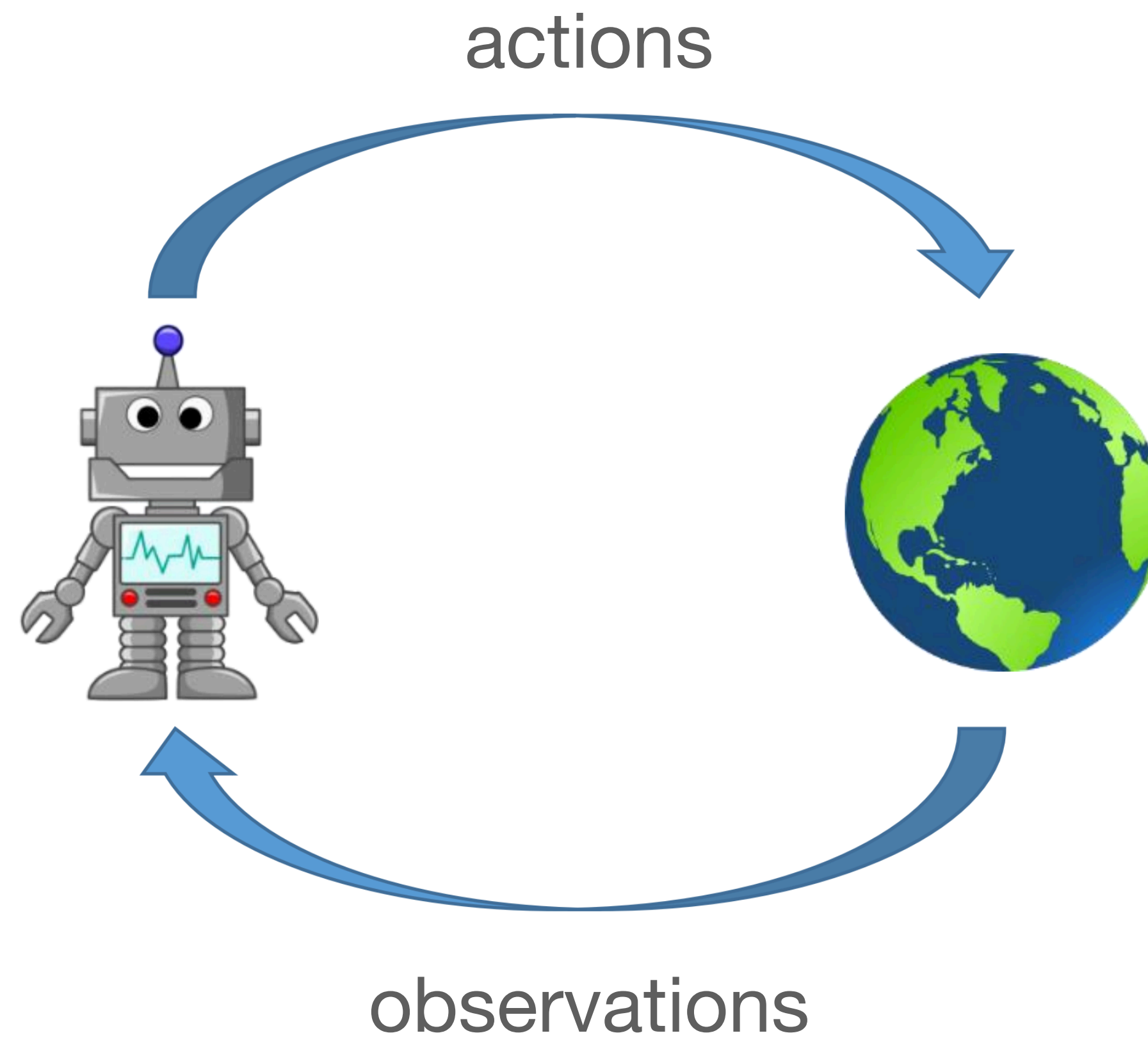
# THE REINFORCEMENT LEARNING LOOP



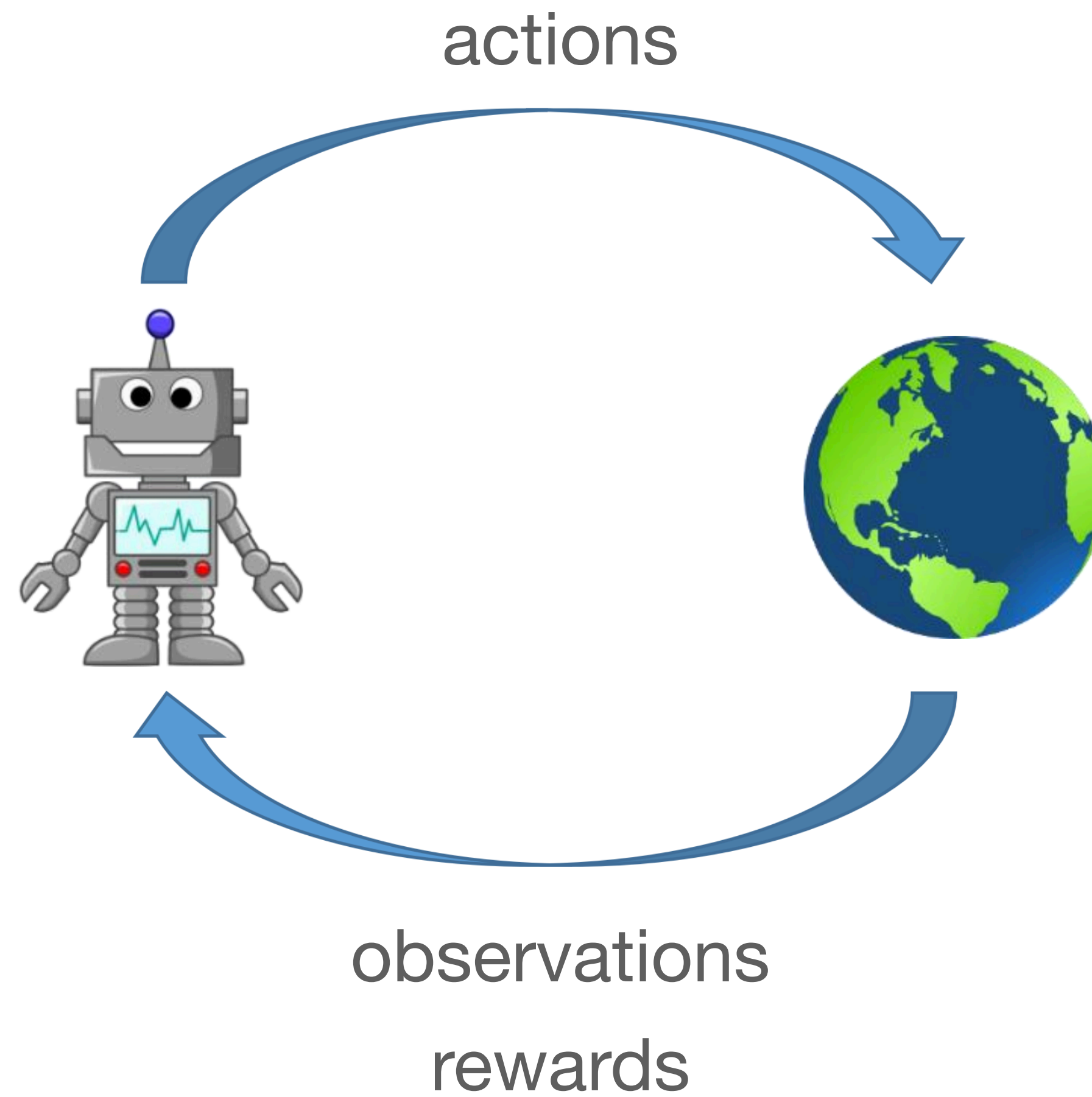
# THE REINFORCEMENT LEARNING LOOP



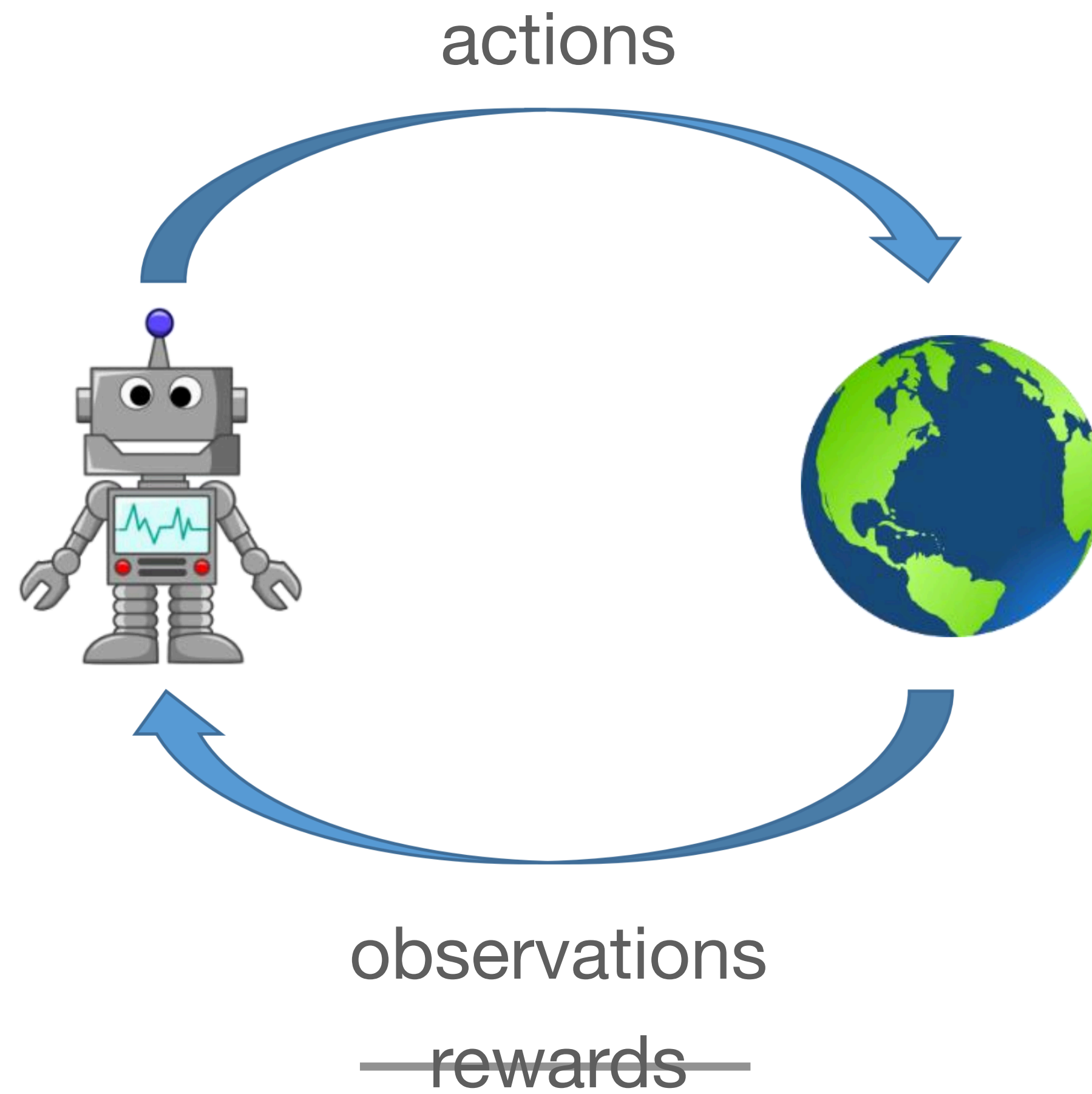
# THE REINFORCEMENT LEARNING LOOP



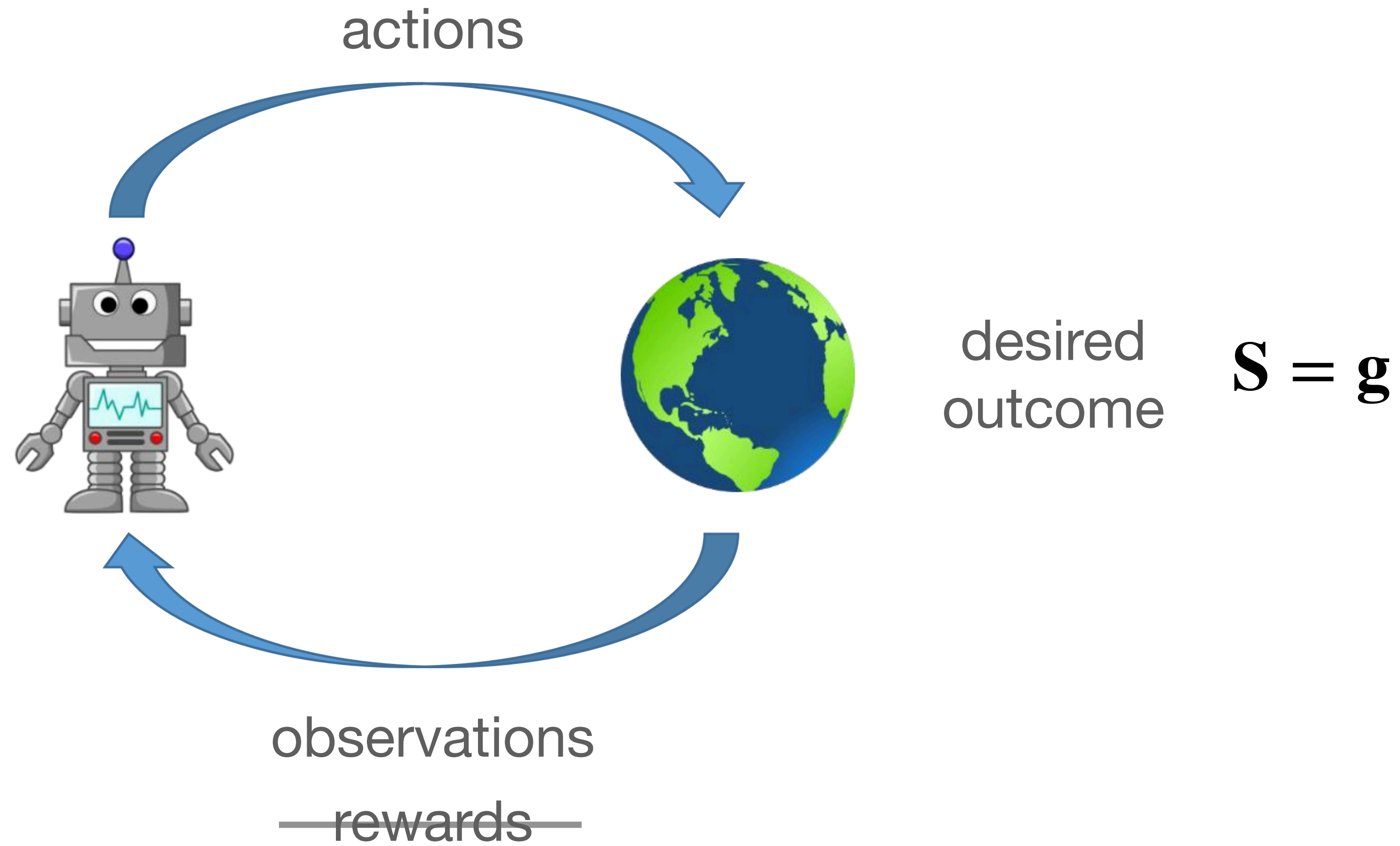
# THE REINFORCEMENT LEARNING LOOP



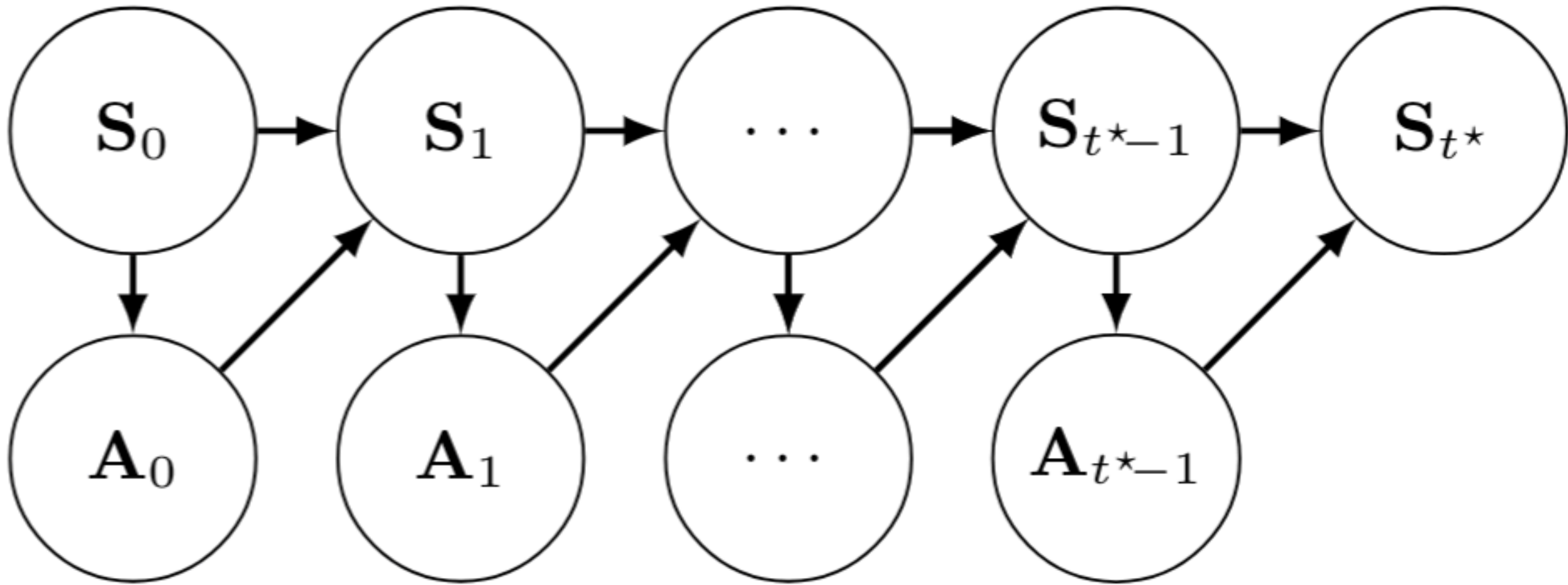
# THE REINFORCEMENT LEARNING LOOP



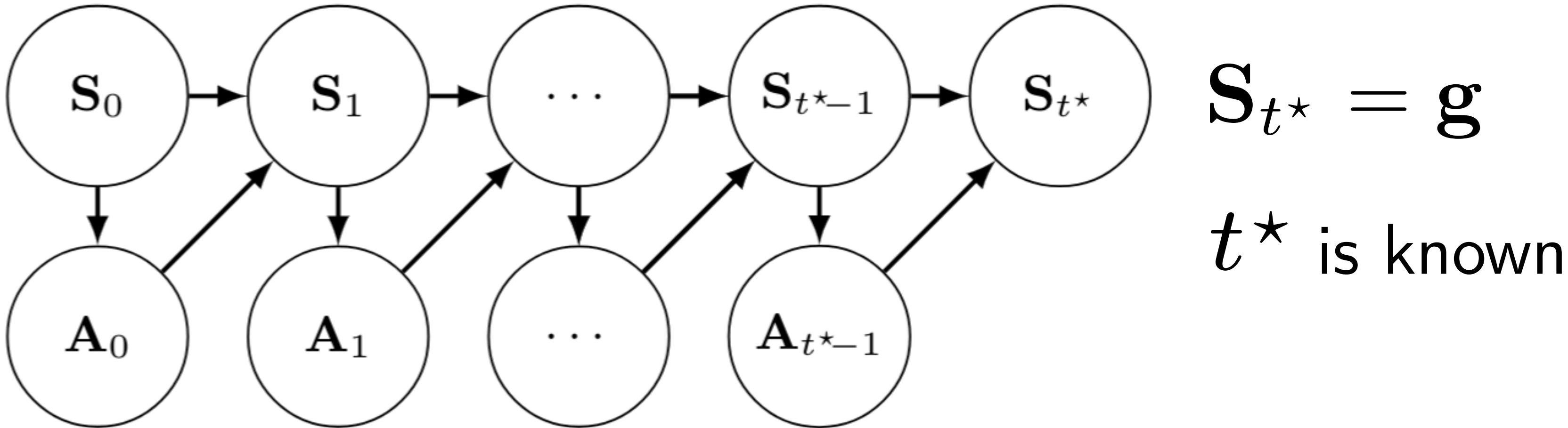
# THE REINFORCEMENT LEARNING LOOP



# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS

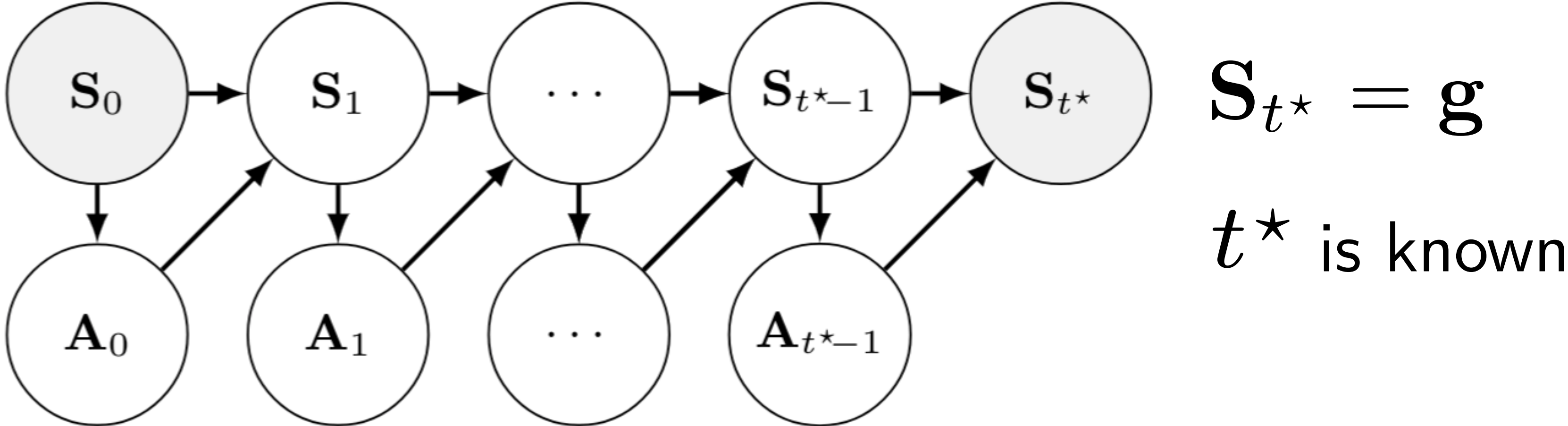


# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS



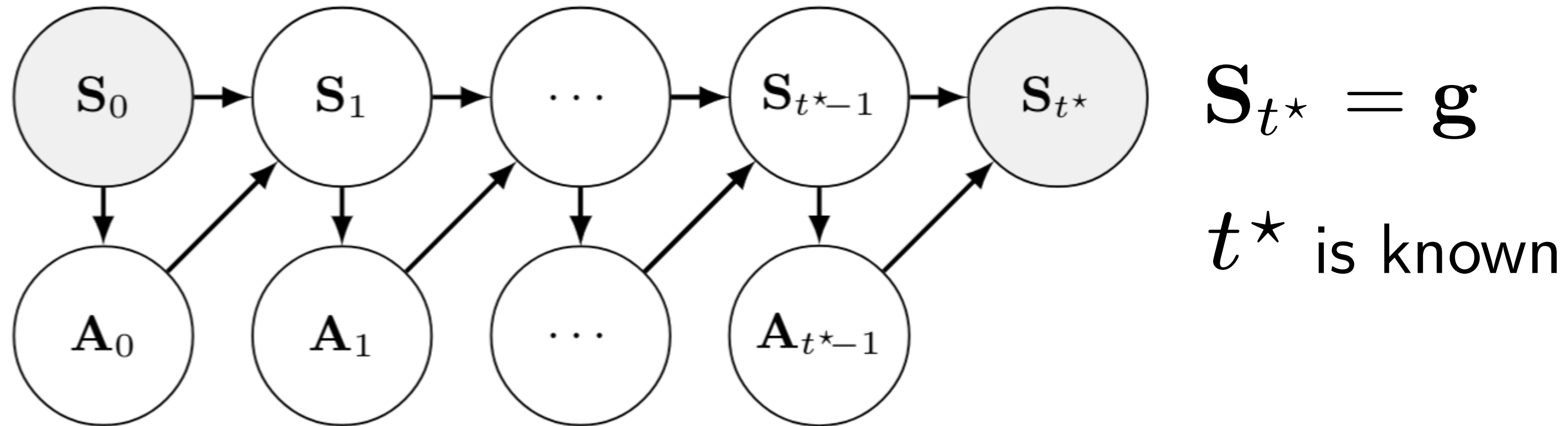


# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS



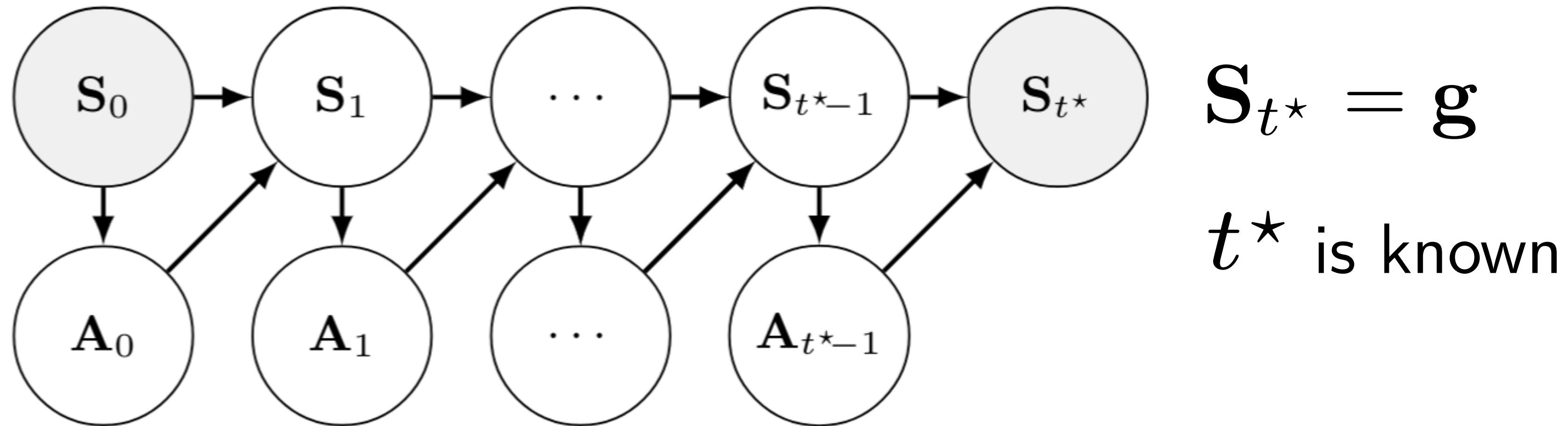
# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS

$$\tilde{\tau}_{0:t} \doteq \{\mathbf{a}_0, \mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_t, \mathbf{a}_t\} \quad t \doteq t^* - 1$$



# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS

$$\tilde{\tau}_{0:t} \doteq \{\mathbf{a}_0, \mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_t, \mathbf{a}_t\} \quad t \doteq t^* - 1$$

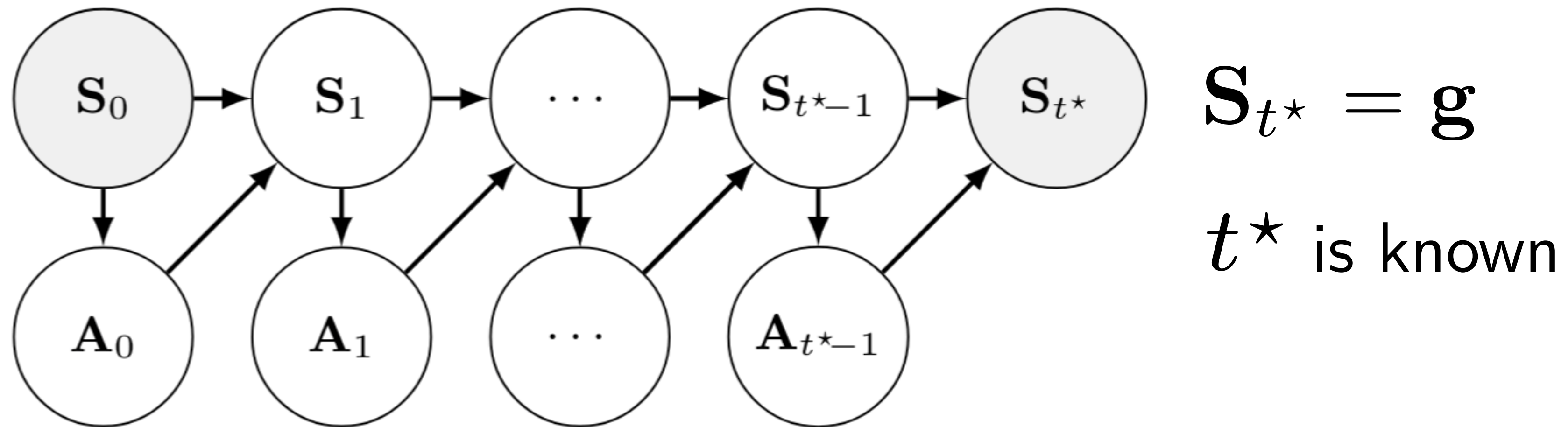


## Goal

- ▶ Find a state-action trajectory distribution leading to  $\mathbf{s}_{t^*} = \mathbf{g}$
- ▶ Infer conditional distribution  $p_{\tilde{\tau}_{0:t} | \mathbf{s}_0, \mathbf{s}_{t^*}}(\cdot | \mathbf{s}_0, \mathbf{g})$

# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS

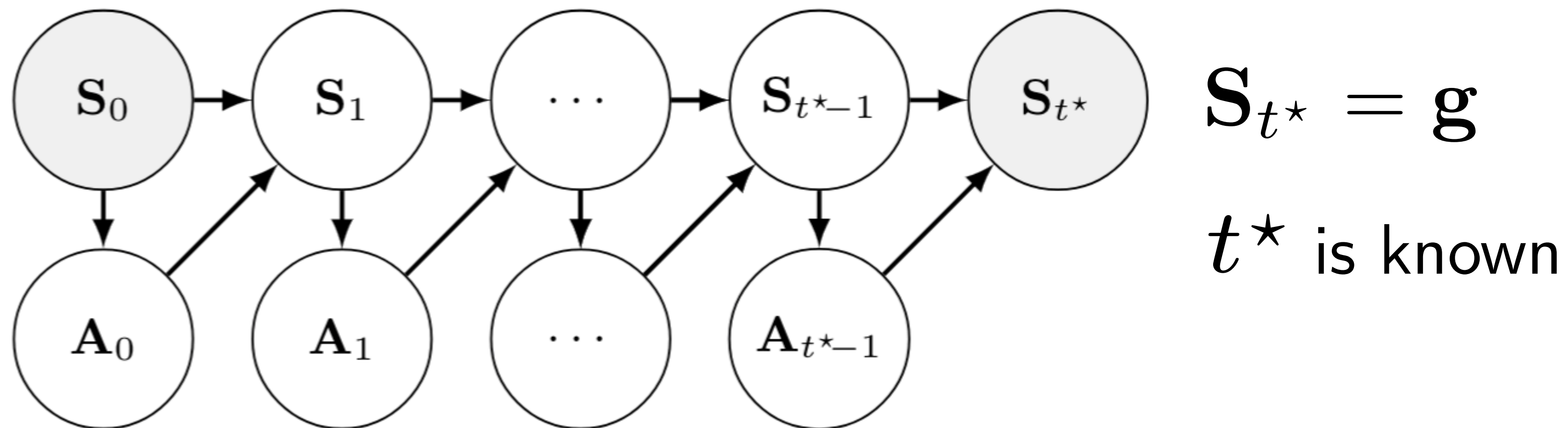
$$\tilde{\tau}_{0:t} \doteq \{\mathbf{a}_0, \mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_t, \mathbf{a}_t\} \quad t \doteq t^* - 1$$



**How?**

# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS

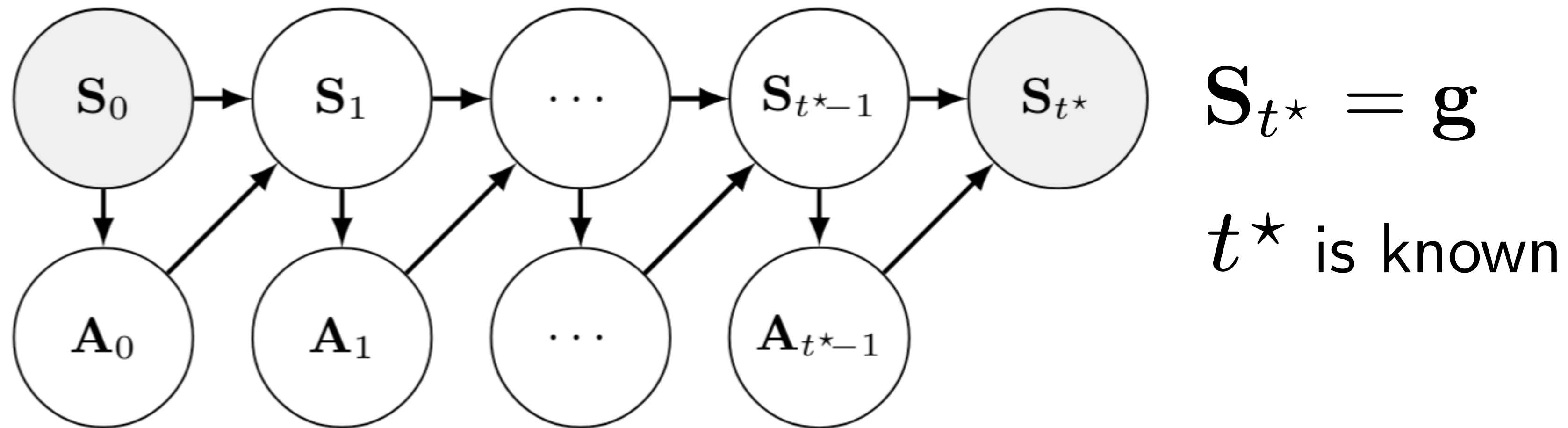
$$\tilde{\tau}_{0:t} \doteq \{\mathbf{a}_0, \mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_t, \mathbf{a}_t\} \quad t \doteq t^* - 1$$



**How?**  $\mathbb{I}\{\mathbf{S}_{t^*} = \mathbf{g}\}$  ?

# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS

$$\tilde{\tau}_{0:t} \doteq \{\mathbf{a}_0, \mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_t, \mathbf{a}_t\} \quad t \doteq t^* - 1$$



**How?**  ~~$\mathbb{I}\{\mathbf{s}_{t^*} = \mathbf{g}\}$~~ ?

- Solve variational problem

$$\min_{q_{\tilde{\tau}_{0:t}|\mathbf{s}_0} \in \hat{\mathcal{Q}}} D_{\text{KL}}(q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\cdot|\mathbf{s}_0) \| p_{\tilde{\tau}_{0:t}|\mathbf{s}_0, \mathbf{s}_{t^*}}(\cdot|\mathbf{s}_0, \mathbf{g}))$$

# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS

**Problem:**

$$\min_{q_{\tilde{\mathcal{T}}_{0:t}|\mathbf{s}_0} \in \hat{\mathcal{Q}}} D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:t}|\mathbf{s}_0}(\cdot|\mathbf{s}_0) \parallel p_{\tilde{\mathcal{T}}_{0:t}|\mathbf{s}_0, \mathbf{s}_{t^*}}(\cdot|\mathbf{s}_0, \mathbf{g}))$$

# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS

## Problem:

$$\min_{q_{\tilde{\mathcal{T}}_{0:t}|\mathbf{s}_0} \in \hat{\mathcal{Q}}} D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:t}|\mathbf{s}_0}(\cdot|\mathbf{s}_0) \parallel p_{\tilde{\mathcal{T}}_{0:t}|\mathbf{s}_0, \mathbf{s}_{t^*}}(\cdot|\mathbf{s}_0, \mathbf{g}))$$

- Intractable!



# INFERRING OUTCOME-DRIVEN TRAJECTORY DISTRIBUTIONS

## Problem:

$$\min_{q_{\tilde{\mathcal{T}}_{0:t}|\mathbf{s}_0} \in \hat{\mathcal{Q}}} D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:t}|\mathbf{s}_0}(\cdot|\mathbf{s}_0) \parallel p_{\tilde{\mathcal{T}}_{0:t}|\mathbf{s}_0, \mathbf{s}_{t^*}}(\cdot|\mathbf{s}_0, \mathbf{g}))$$

- Intractable!

## Instead:

- Derive an **equivalent tractable** variational objective

# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE (FINITE HORIZON)

**Deriving a tractable variational lower bound**

# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE (FINITE HORIZON)

## Deriving a tractable variational lower bound

- Define variational distribution

$$q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\tilde{\tau}_{0:t} | \mathbf{s}_0) \doteq \pi(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$

# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE (FINITE HORIZON)

## Deriving a tractable variational lower bound

- ▶ Define variational distribution

$$q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\tilde{\tau}_{0:t} | \mathbf{s}_0) \doteq \pi(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$

- ▶ Then:

$$D_{\text{KL}}(q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\tau}_{0:t}|\mathbf{s}_0, \mathbf{s}_{t^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = -\bar{\mathcal{F}}(\pi, \mathbf{s}_0, \mathbf{g}) + \log p(\mathbf{g}|\mathbf{s}_0)$$

# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE (FINITE HORIZON)

## Deriving a tractable variational lower bound

- ▶ Define variational distribution

$$q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\tilde{\tau}_{0:t} | \mathbf{s}_0) \doteq \pi(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$

- ▶ Then:

$$D_{\text{KL}}(q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\tau}_{0:t}|\mathbf{s}_0, \mathbf{s}_{t^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = -\bar{\mathcal{F}}(\pi, \mathbf{s}_0, \mathbf{g}) + \log p(\mathbf{g}|\mathbf{s}_0)$$

with

$$\bar{\mathcal{F}}(\pi, \mathbf{s}_0, \mathbf{g}) \doteq \mathbb{E}_{q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\tilde{\tau}_{0:t} | \mathbf{s}_0)} \left[ \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t'=0}^{t-1} D_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \| p(\cdot | \mathbf{s}_{t'})) \right]$$

# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE (FINITE HORIZON)

## Deriving a tractable variational lower bound

- ▶ Define variational distribution

$$q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\tilde{\tau}_{0:t} | \mathbf{s}_0) \doteq \pi(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$

- ▶ **Hence:**

$$\arg \min_{\pi \in \Pi} D_{\text{KL}}(q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\cdot | \mathbf{s}_0) || p_{\tilde{\tau}_{0:t}|\mathbf{s}_0, \mathbf{s}_{t^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = \arg \max_{\pi \in \Pi} \bar{\mathcal{F}}(\pi, \mathbf{s}_0, \mathbf{g}).$$

with

$$\bar{\mathcal{F}}(\pi, \mathbf{s}_0, \mathbf{g}) \doteq \mathbb{E}_{q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}(\tilde{\tau}_{0:t} | \mathbf{s}_0)} \left[ \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t'=0}^{t-1} D_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) || p(\cdot | \mathbf{s}_{t'})) \right]$$

# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE (FINITE HORIZON)

**What reward function is this objective optimizing?**

- Variational objective

$$\bar{\mathcal{F}}(\pi, \mathbf{s}_0, \mathbf{g}) \doteq \mathbb{E}_{q_{\tilde{\tau}_{0:t}|\mathbf{s}_0}}(\tilde{\tau}_{0:t} | \mathbf{s}_0) \left[ \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t'=0}^{t-1} D_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) || p(\cdot | \mathbf{s}_{t'})) \right]$$

# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE (FINITE HORIZON)

## What reward function is this objective optimizing?

- Variational objective

$$\bar{\mathcal{F}}(\pi, \mathbf{s}_0, \mathbf{g}) \doteq \mathbb{E}_{q_{\tilde{\tau}_{0:t} | \mathbf{s}_0}}(\tilde{\tau}_{0:t} | \mathbf{s}_0) \left[ \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t'=0}^{t-1} D_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) || p(\cdot | \mathbf{s}_{t'})) \right]$$

- Equivalent to KL-regularized RL with rewards

$$r(\mathbf{s}_{t'}, \mathbf{a}_{t'}, \mathbf{g}, t') \doteq \mathbb{I}\{t' = t\} \log p_d(\mathbf{g} | \mathbf{s}_{t'}, \mathbf{a}_{t'})$$



# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE (FINITE HORIZON)

## What reward function is this objective optimizing?

- Variational objective

$$\bar{\mathcal{F}}(\pi, \mathbf{s}_0, \mathbf{g}) \doteq \mathbb{E}_{q_{\tilde{\tau}_{0:t} | \mathbf{s}_0}}(\tilde{\tau}_{0:t} | \mathbf{s}_0) \left[ \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) + \sum_{t'=0}^{t-1} D_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) || p(\cdot | \mathbf{s}_{t'})) \right]$$

- Equivalent to KL-regularized RL with rewards

$$r(\mathbf{s}_{t'}, \mathbf{a}_{t'}, \mathbf{g}, t') \doteq \mathbb{I}\{t' = t\} \log p_d(\mathbf{g} | \mathbf{s}_{t'}, \mathbf{a}_{t'})$$

- Sparse rewards given by transition dynamics

# INFINITE-HORIZON OUTCOME-DRIVEN MDP

**But: We don't care *when* the outcome is achieved**

$$p_{\tilde{\mathcal{T}}_{0:T}, \mathbf{s}_{T^*}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, \mathbf{g}, t | \mathbf{s}_0) = p_T(t) p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) p(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$

# INFINITE-HORIZON OUTCOME-DRIVEN MDP

**But: We don't care *when* the outcome is achieved**

- ▶ Treat **termination time** as a **random variable**:

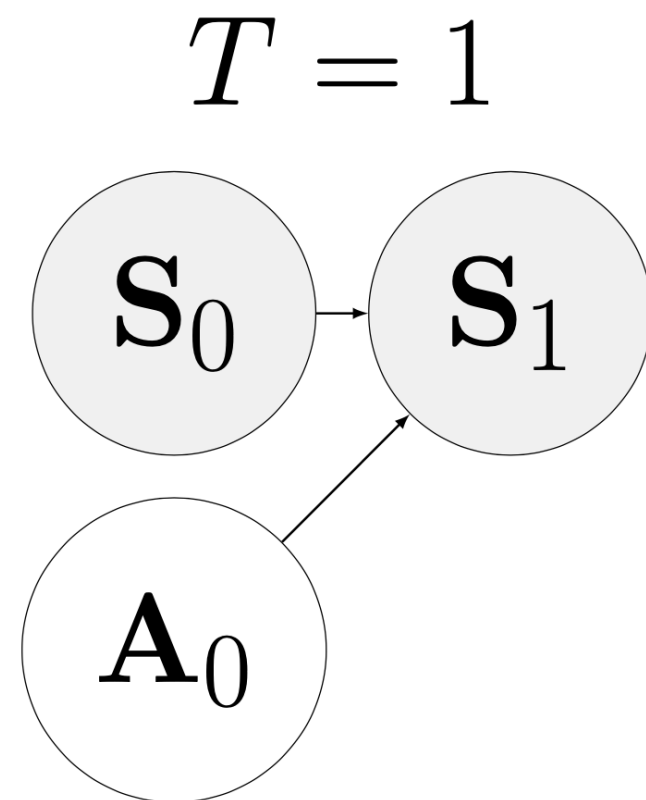
$$p_{\tilde{\mathcal{T}}_{0:T}, \mathbf{s}_{T^*}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, \mathbf{g}, t | \mathbf{s}_0) = p_T(t) p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) p(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$

# INFINITE-HORIZON OUTCOME-DRIVEN MDP

**But: We don't care *when* the outcome is achieved**

- ▶ Treat **termination time** as a **random variable**:

$$p_{\tilde{\tau}_{0:T}, \mathbf{s}_{T^*}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, \mathbf{g}, t | \mathbf{s}_0) = p_T(t) p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) p(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$

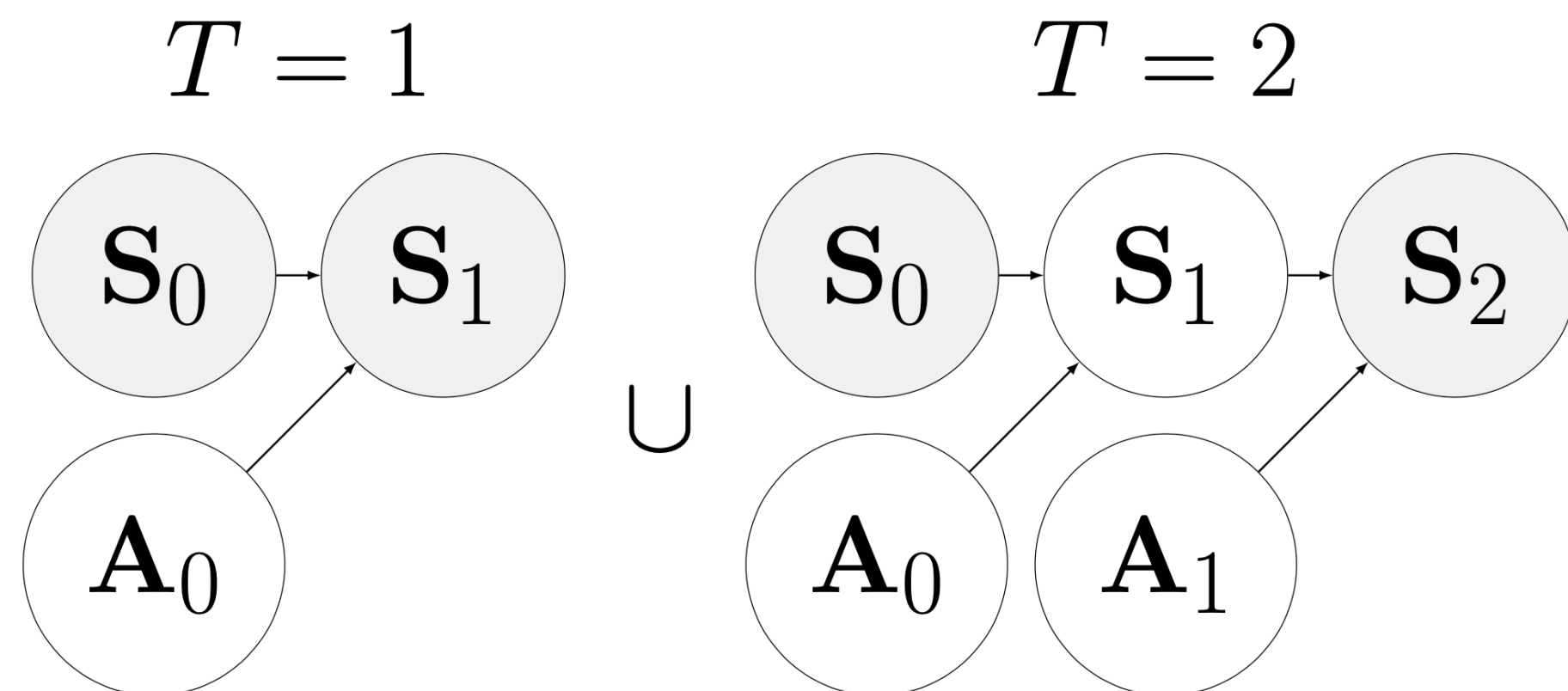


# INFINITE-HORIZON OUTCOME-DRIVEN MDP

**But:** We don't care *when* the outcome is achieved

- ▶ Treat **termination time** as a **random variable**:

$$p_{\tilde{\tau}_{0:T}, \mathbf{s}_{T^*}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, \mathbf{g}, t | \mathbf{s}_0) = p_T(t) p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) p(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$

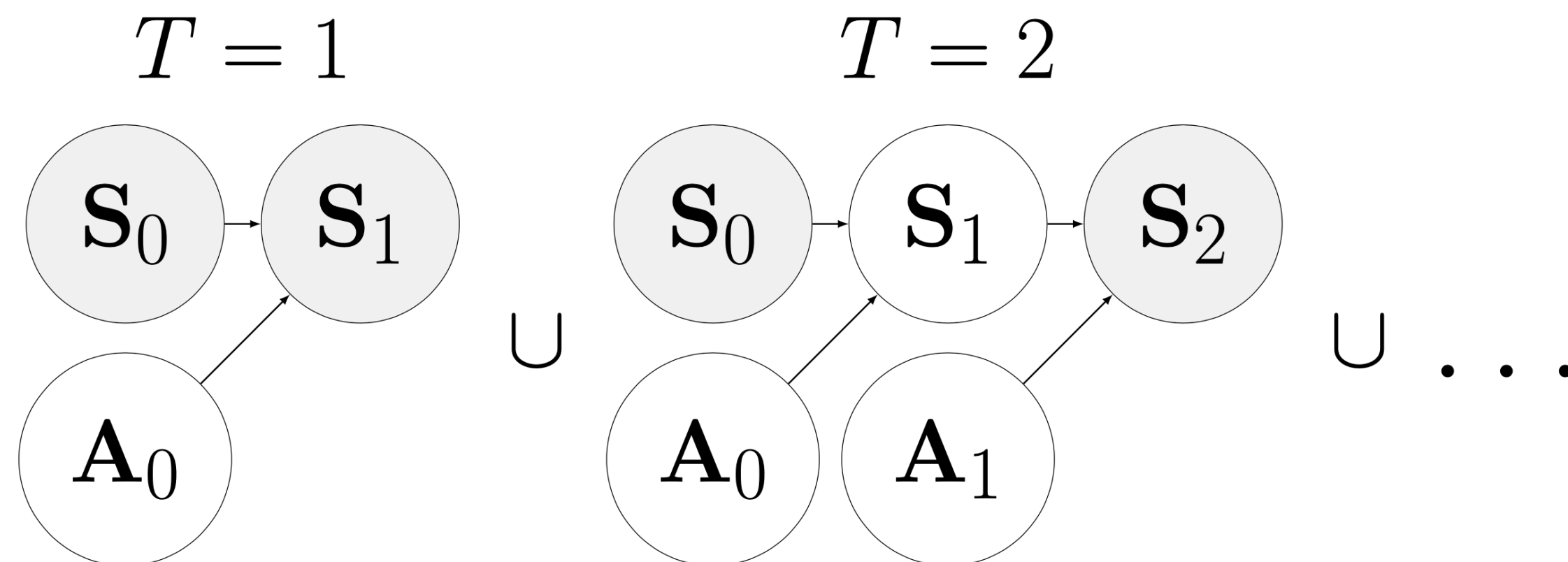


# INFINITE-HORIZON OUTCOME-DRIVEN MDP

**But:** We don't care *when* the outcome is achieved

- ▶ Treat **termination time** as a **random variable**:

$$p_{\tilde{\tau}_{0:T}, \mathbf{s}_{T^*}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, \mathbf{g}, t | \mathbf{s}_0) = p_T(t) p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) p(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$



# INFERENCE IN AN INFINITE-HORIZON OUTCOME-DRIVEN MDP

**But:** We don't care *when* the outcome is achieved

- ▶ Treat **termination time** as a **random variable**:

$$p_{\tilde{\mathcal{T}}_{0:T}, \mathbf{s}_{T^*}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, \mathbf{g}, t | \mathbf{s}_0) = p_T(t) p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) p(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} | \mathbf{s}_{t'})$$

- ▶ Perform inference over **states**, **actions**, and the **termination time**

$$\min_{q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0} \in \mathcal{Q}} D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g}))$$

## **Infinite-horizon variational objective**

- Define variational distribution



## Infinite-horizon variational objective

- Define variational distribution

$$q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, t | \mathbf{s}_0) = q_{\tilde{\mathcal{T}}_{0:T} | T, \mathbf{s}_0}(\tilde{\tau}_{0:t} | t, \mathbf{s}_0) q_T(t)$$

# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE

## Infinite-horizon variational objective

- Define variational distribution

$$q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, t | \mathbf{s}_0) = q_{\tilde{\mathcal{T}}_{0:T} | T, \mathbf{s}_0}(\tilde{\tau}_{0:t} | t, \mathbf{s}_0) q_T(t)$$

- Then:

$$D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = \log p(\mathbf{g} | \mathbf{s}_0) - \mathcal{F}(\pi, q_T, \mathbf{s}_0, \mathbf{g})$$

with

$$\mathcal{F}(\pi, q_T, \mathbf{s}_0, \mathbf{g})$$

$$\doteq \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\mathcal{T}}_{0:T} | T, \mathbf{s}_0}(\tilde{\tau}_{0:t} | t, \mathbf{s}_0)} \left[ \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) - D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0)) \right]$$

# OUTCOME-DRIVEN VARIATIONAL OBJECTIVE

## Infinite-horizon variational objective

- Define variational distribution

$$q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, t | \mathbf{s}_0) = q_{\tilde{\tau}_{0:T} | T, \mathbf{s}_0}(\tilde{\tau}_{0:t} | t, \mathbf{s}_0) q_T(t)$$

- Then:

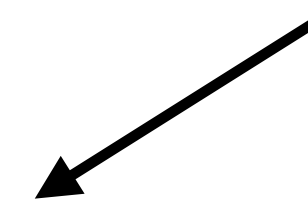
$$D_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = \log p(\mathbf{g} | \mathbf{s}_0) - \mathcal{F}(\pi, q_T, \mathbf{s}_0, \mathbf{g})$$

with

$$\mathcal{F}(\pi, q_T, \mathbf{s}_0, \mathbf{g})$$

$$\doteq \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T} | T, \mathbf{s}_0}(\tilde{\tau}_{0:t} | t, \mathbf{s}_0)} \left[ \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) - D_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0)) \right]$$

not factorized



# OUTCOME-DRIVEN VARIATIONAL DISTRIBUTION

## Infinite-horizon variational objective

- ▶ Define variational distribution

$$q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, t | \mathbf{s}_0) = q_{\tilde{\mathcal{T}}_{0:T} | T, \mathbf{s}_0}(\tilde{\tau}_{0:t} | t, \mathbf{s}_0) q_T(t)$$

$$q_T(t) = q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0)$$

- ▶ Goal: Derive recursive variational objective

# OUTCOME-DRIVEN VARIATIONAL INFERENCE AS TD LEARNING

**Theorem 1.** (informal)

# OUTCOME-DRIVEN VARIATIONAL INFERENCE AS TD LEARNING

## Theorem 1. (informal)

- Define variational distributions  $q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, t | \mathbf{s}_0)$  and  $q_T(t)$

# OUTCOME-DRIVEN VARIATIONAL INFERENCE AS TD LEARNING

## Theorem 1. (informal)

- ▶ Define variational distributions  $q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, t | \mathbf{s}_0)$  and  $q_T(t)$
- ▶ Then

$$V^\pi(\mathbf{s}_t, \mathbf{g}; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T)] - D_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t))$$

# OUTCOME-DRIVEN VARIATIONAL INFERENCE AS TD LEARNING

## Theorem 1. (informal)

▸ Define variational distributions  $q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, t | \mathbf{s}_0)$  and  $q_T(t)$

▸ Then

$$V^\pi(\mathbf{s}_t, \mathbf{g}; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T)] - D_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t))$$

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T) \doteq r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1}, \mathbf{g}; \pi, q_T)]$$



# OUTCOME-DRIVEN VARIATIONAL INFERENCE AS TD LEARNING

## Theorem 1. (informal)

▸ Define variational distributions  $q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, t | \mathbf{s}_0)$  and  $q_T(t)$

▸ Then

$$V^\pi(\mathbf{s}_t, \mathbf{g}; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T)] - D_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t))$$

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T) \doteq r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1}, \mathbf{g}; \pi, q_T)]$$

$$r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) \doteq q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) - D_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}})$$

# OUTCOME-DRIVEN VARIATIONAL INFERENCE AS TD LEARNING

## Theorem 1. (informal)

▸ Define variational distributions  $q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\tilde{\tau}_{0:t}, t | \mathbf{s}_0)$  and  $q_T(t)$

▸ Then

$$V^\pi(\mathbf{s}_t, \mathbf{g}; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T)] - D_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t))$$

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T) \doteq r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1}, \mathbf{g}; \pi, q_T)]$$

$$r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) \doteq q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) - D_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}})$$

with

$$D_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = -V^\pi(\mathbf{s}_0, \mathbf{g}; q_T) + \log p(\mathbf{g} | \mathbf{s}_0)$$

# A LOWER BOUND ON THE LOG MARGINAL LIKELIHOOD

**What does this mean?**

# A LOWER BOUND ON THE LOG MARGINAL LIKELIHOOD

## What does this mean?

▸ Since

$$D_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = -V^\pi(\mathbf{s}_0, \mathbf{g}; q_T) + \log p(\mathbf{g} | \mathbf{s}_0),$$

the **value function** is a **variational lower bound** on  $\log p(\mathbf{g} | \mathbf{s}_0)$

# A LOWER BOUND ON THE LOG MARGINAL LIKELIHOOD

## What does this mean?

▸ Since

$$D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = -V^\pi(\mathbf{s}_0, \mathbf{g}; q_T) + \log p(\mathbf{g} | \mathbf{s}_0),$$

the **value function** is a **variational lower bound** on  $\log p(\mathbf{g} | \mathbf{s}_0)$  and

$$\arg \min_{\pi \in \Pi, q_T \in \mathcal{Q}_T} \{D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g}))\}$$

# A LOWER BOUND ON THE LOG MARGINAL LIKELIHOOD

## What does this mean?

► Since

$$D_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = -V^\pi(\mathbf{s}_0, \mathbf{g}; q_T) + \log p(\mathbf{g} | \mathbf{s}_0),$$

the **value function** is a **variational lower bound** on  $\log p(\mathbf{g} | \mathbf{s}_0)$  and

$$\begin{aligned} & \arg \min_{\pi \in \Pi, q_T \in \mathcal{Q}_T} \{D_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g}))\} \\ &= \arg \max_{\pi \in \Pi, q_T \in \mathcal{Q}_T} \mathcal{F}(\pi, q_T, \mathbf{s}_0, \mathbf{g}) \end{aligned}$$

# A LOWER BOUND ON THE LOG MARGINAL LIKELIHOOD

## What does this mean?

► Since

$$D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g})) = -V^\pi(\mathbf{s}_0, \mathbf{g}; q_T) + \log p(\mathbf{g} | \mathbf{s}_0),$$

the **value function** is a **variational lower bound** on  $\log p(\mathbf{g} | \mathbf{s}_0)$  and

$$\begin{aligned} & \arg \min_{\pi \in \Pi, q_T \in \mathcal{Q}_T} \{D_{\text{KL}}(q_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\mathcal{T}}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \mathbf{s}_0, \mathbf{g}))\} \\ &= \arg \max_{\pi \in \Pi, q_T \in \mathcal{Q}_T} \mathcal{F}(\pi, q_T, \mathbf{s}_0, \mathbf{g}) \\ &= \arg \max_{\pi \in \Pi, q_T \in \mathcal{Q}_T} V^\pi(\mathbf{s}_0, \mathbf{g}; q_T). \end{aligned}$$

# OUTCOME-DRIVEN BELLMAN OPERATOR

## Outcome-Driven Bellman Operator

$$\mathcal{T}^\pi Q(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T) \doteq r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V(\mathbf{s}_{t+1}, \mathbf{g}; q_T)]$$

with

$$V(\mathbf{s}_t, \mathbf{g}; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T)] + D_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t))$$

$$r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) \doteq q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) - D_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}})$$



# OUTCOME-DRIVEN BELLMAN OPERATOR

## Outcome-Driven Bellman Operator

$$\mathcal{T}^\pi Q(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T) \doteq r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V(\mathbf{s}_{t+1}, \mathbf{g}; q_T)]$$

with

$$V(\mathbf{s}_t, \mathbf{g}; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_T)] + D_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t))$$

$$r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) \doteq q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) - D_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}})$$

- ▶ **Dense, shaped reward function + dynamic discount factor**

# OUTCOME-DRIVEN POLICY ITERATION

## **Theorem 2.** (informal)

- ▶ Assume that the MDP is ergodic and set of available action is **finite**.

# OUTCOME-DRIVEN POLICY ITERATION

## **Theorem 2.** (informal)

- ▶ Assume that the MDP is ergodic and set of available action is **finite**.
- ▶ Consider the **outcome-driven Bellman operator**.

# OUTCOME-DRIVEN POLICY ITERATION

## Theorem 2. (informal)

- ▶ Assume that the MDP is ergodic and set of available action is **finite**.
- ▶ Consider the **outcome-driven Bellman operator**.
- ▶ Consider the **optimal variational distribution** over the termination time

# OUTCOME-DRIVEN POLICY ITERATION

## Theorem 2. (informal)

- ▶ Assume that the MDP is ergodic and set of available action is **finite**.
- ▶ Consider the **outcome-driven Bellman operator**.
- ▶ Consider the **optimal variational distribution** over the termination time
- ▶ Alternating between **policy evaluation** and **policy improvement** will result in an **optimal policy**

# OUTCOME-DRIVEN ACTOR CRITIC

---

**Algorithm 1** ODAC: Outcome-Driven Actor–Critic

---

- 1: Initialize policy  $\pi_\theta$ , replay buffer  $\mathcal{R}$ ,  $Q$ -function  $Q_\phi$ , and dynamics model  $p_\psi$ .
  - 2: **for** iteration  $i = 1, 2, \dots$  **do**
  - 3:     Collect on-policy samples to add to  $\mathcal{R}$  by sampling  $\mathbf{g}$  from environment and executing  $\pi$ .
  - 4:     Sample batch  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', \mathbf{g})$  from  $\mathcal{R}$ .
  - 5:     Compute approximate reward and optimal weights
  - 6:     Update  $Q_\phi$ ,  $\pi_\theta$ , and  $p_\psi$
  - 7: **end for**
- 

with  $\hat{r}(\mathbf{s}_t, \mathbf{a}_t, \mathbf{g}; q_\Delta) \doteq \hat{q}_{\Delta_{t+1}}(\Delta_{t+1} = 1) \log p_\psi(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) - D_{\text{KL}}(q_{\Delta_t} \| p_{\Delta_t})$

$$\mathcal{F}_Q(\phi) = \mathbb{E} \left[ \left( Q_\phi(\mathbf{s}, \mathbf{a}, \mathbf{g}) - (\hat{r}(\mathbf{s}, \mathbf{a}, \mathbf{g}; q_\Delta) + q_{\Delta_t}(\Delta_t = 0) \hat{V}(\mathbf{s}', \mathbf{g})) \right)^2 \right]$$

$$\mathcal{F}_\pi(\theta) = -\mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \pi_\theta(\cdot | \mathbf{s}; \mathbf{g})} [Q_\phi(\mathbf{s}, \mathbf{a}, \mathbf{g}) - \log \pi_\theta(\mathbf{a} | \mathbf{s}; \mathbf{g})]$$

$$\mathcal{F}_p(\psi) = \mathbb{E}_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \sim \mathcal{D}} [\log p_\psi(\mathbf{s}' | \mathbf{s}, \mathbf{a})]$$

# OUTCOME-DRIVEN REWARD ESTIMATION

## Estimation of transition dynamics

- ▶ Objective is defined in terms of **likelihood of achieving the outcome**

## Estimation of transition dynamics

- ▶ Objective is defined in terms of **likelihood of achieving the outcome**
- ▶ Need to know transition dynamics



## Estimation of transition dynamics

- ▶ Objective is defined in terms of **likelihood of achieving the outcome**
- ▶ Need to know transition dynamics
- ▶ If transition dynamics are unknown, estimation is needed:

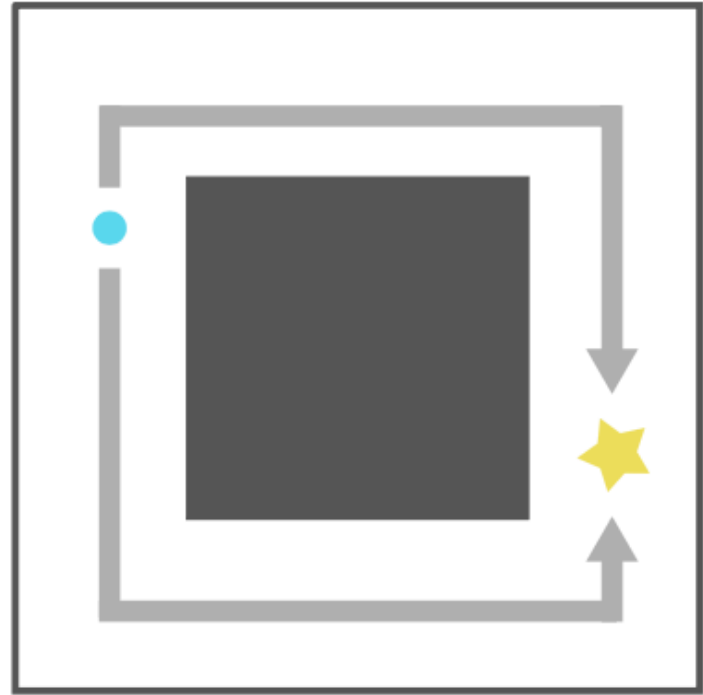
## Estimation of transition dynamics

- ▶ Objective is defined in terms of **likelihood of achieving the outcome**
- ▶ Need to know transition dynamics
- ▶ If transition dynamics are unknown, estimation is needed:
  - ▶ Fix parametric model

## Estimation of transition dynamics

- ▶ Objective is defined in terms of **likelihood of achieving the outcome**
- ▶ Need to know transition dynamics
- ▶ If transition dynamics are unknown, estimation is needed:
  - ▶ Fix parametric model
  - ▶ Learn parametric model via  $\mathcal{F}_p(\psi) = \mathbb{E}_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \sim \mathcal{D}}[\log p_\psi(\mathbf{s}' | \mathbf{s}, \mathbf{a})]$

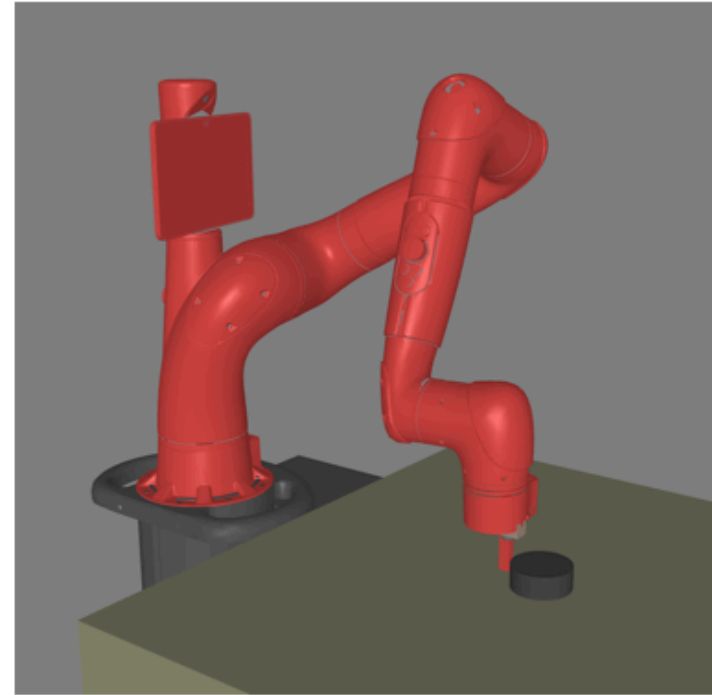
# EXPERIMENTS



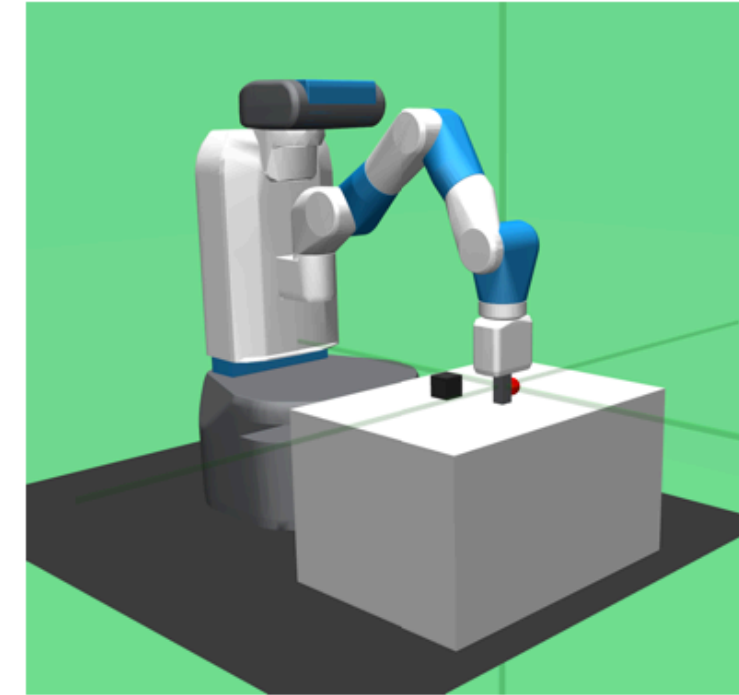
(a) Box 2D



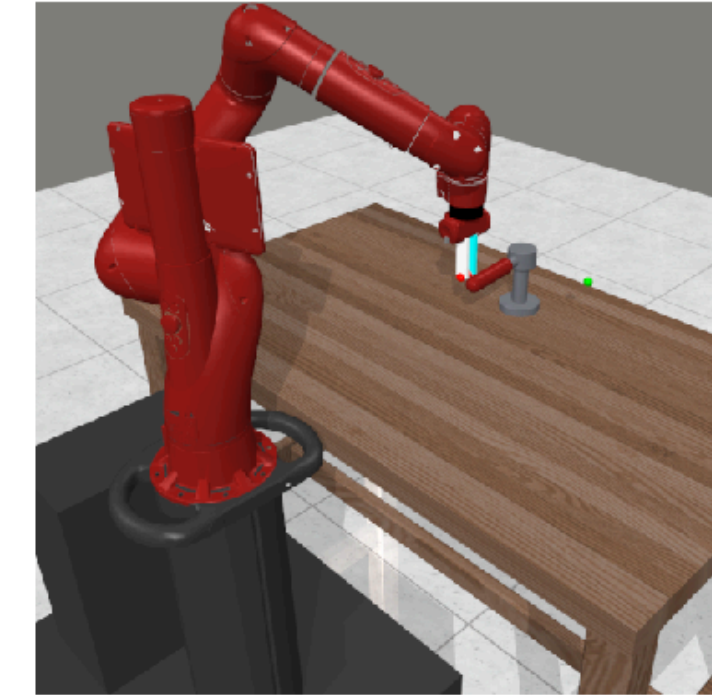
(b) Ant



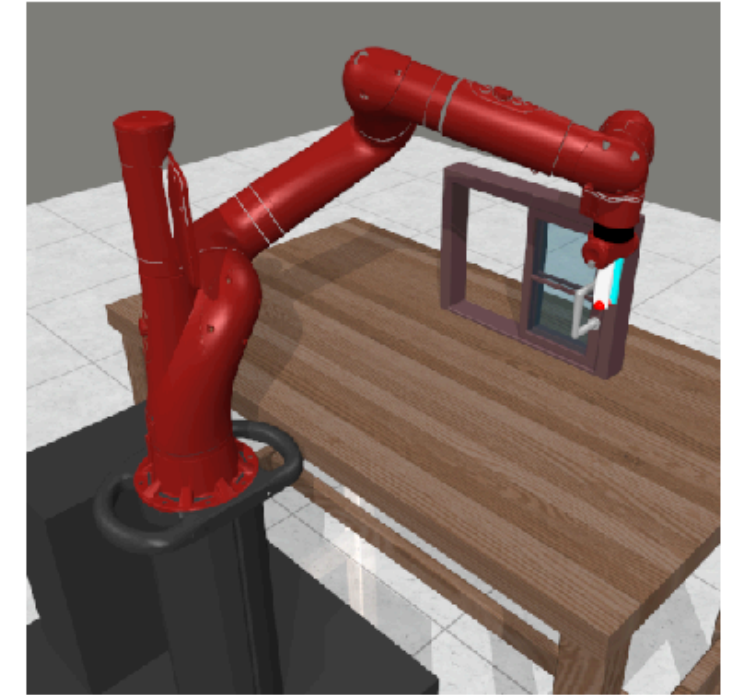
(c) Sawyer Push



(d) Fetch Push



(e) Faucet

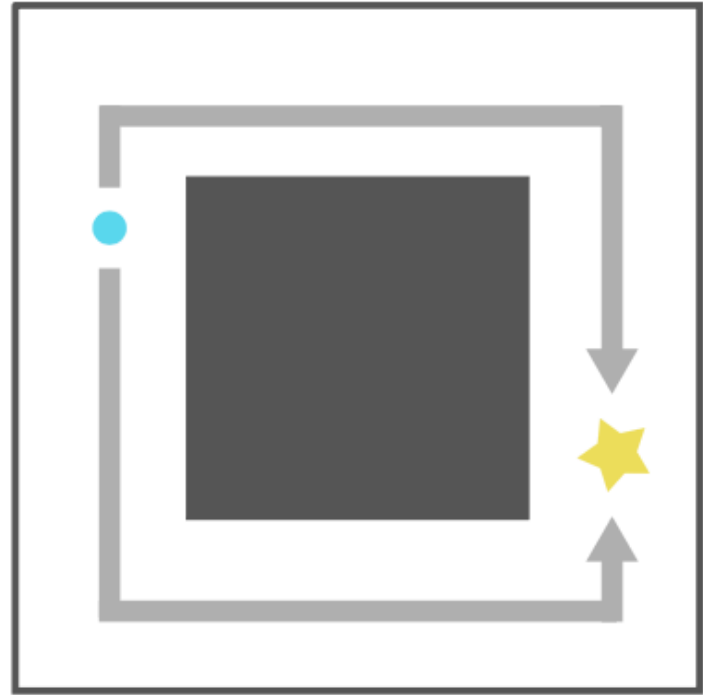


(f) Window

## Setting

- ▶ No oracle goal sampling

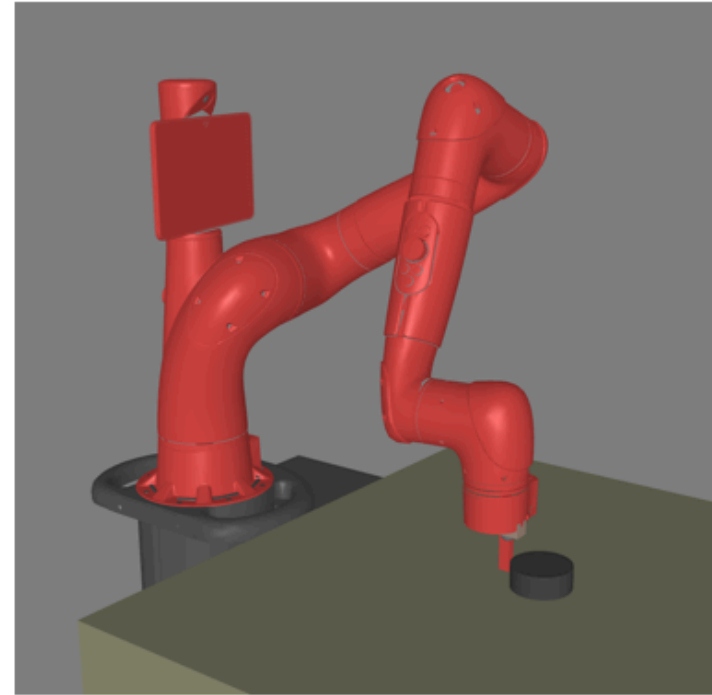
# EXPERIMENTS



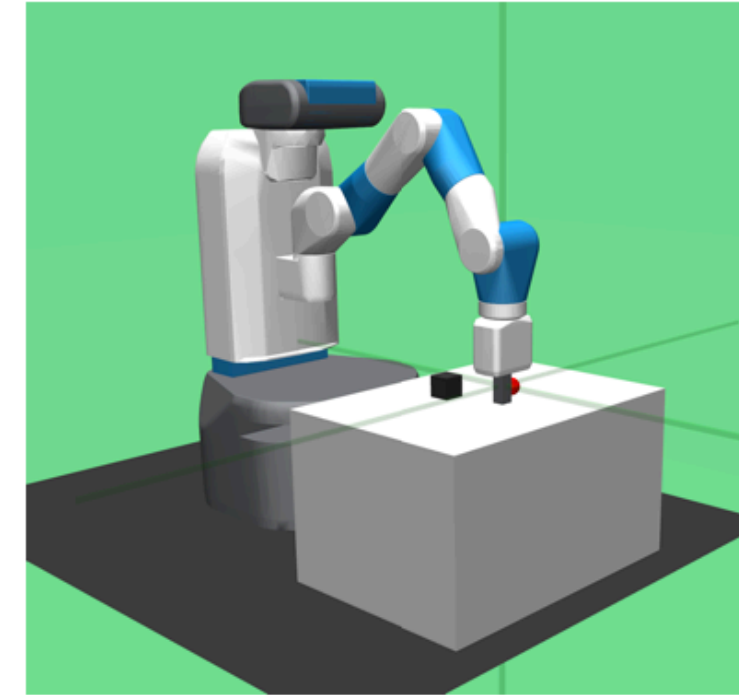
(a) Box 2D



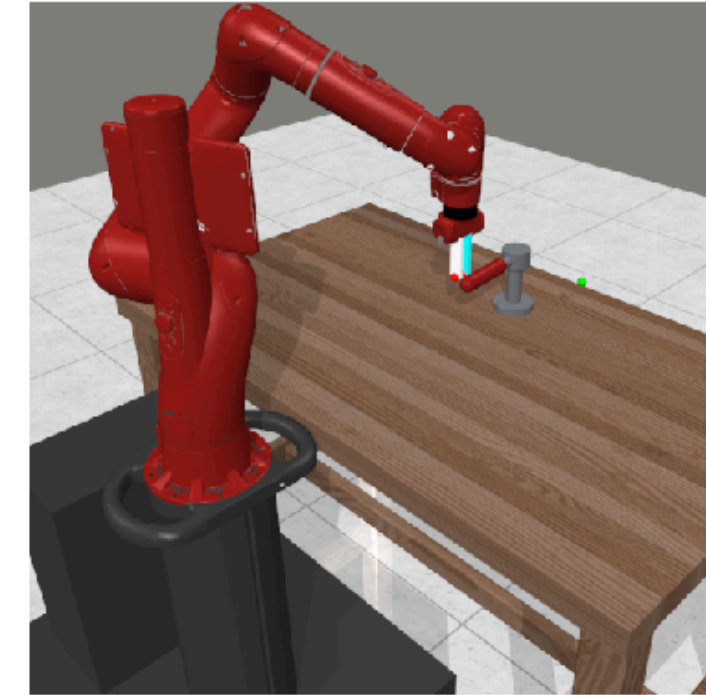
(b) Ant



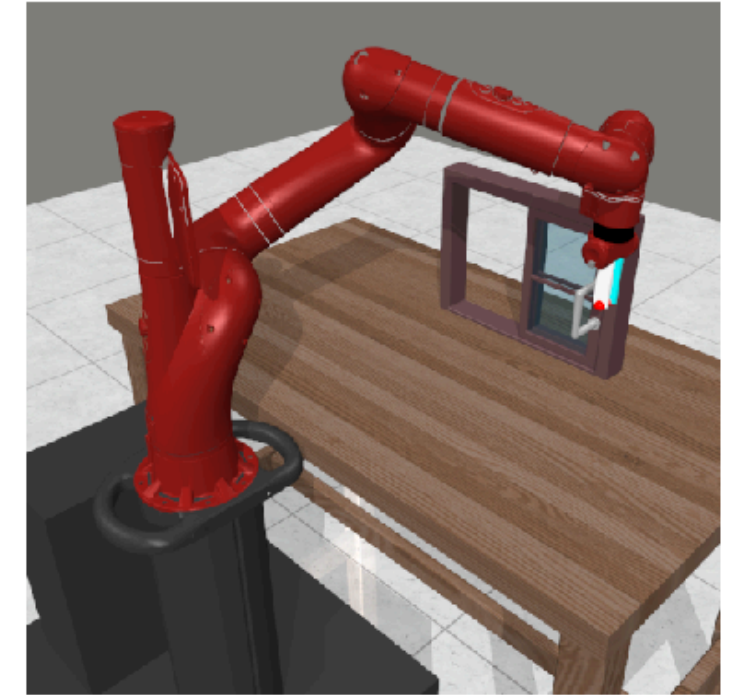
(c) Sawyer Push



(d) Fetch Push



(e) Faucet

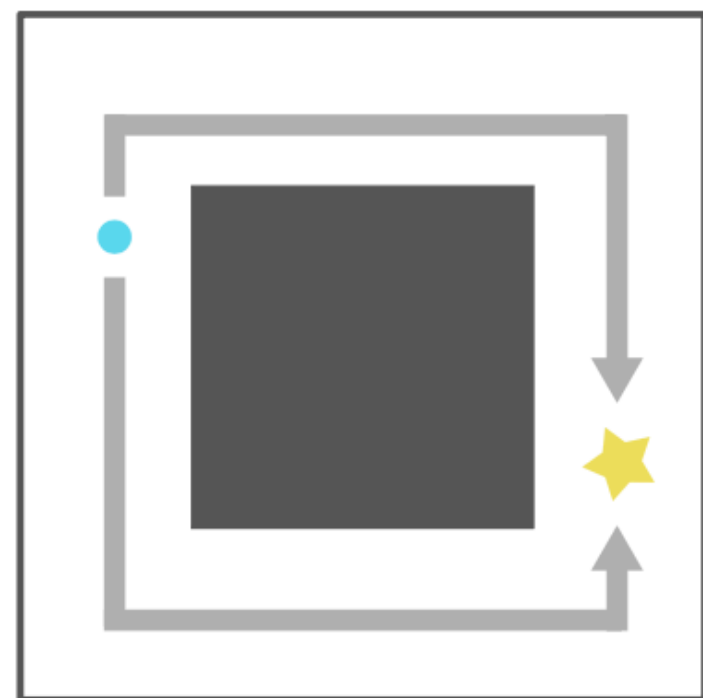


(f) Window

## Setting

- ▶ No oracle goal sampling
- ▶ Future-style relabeling

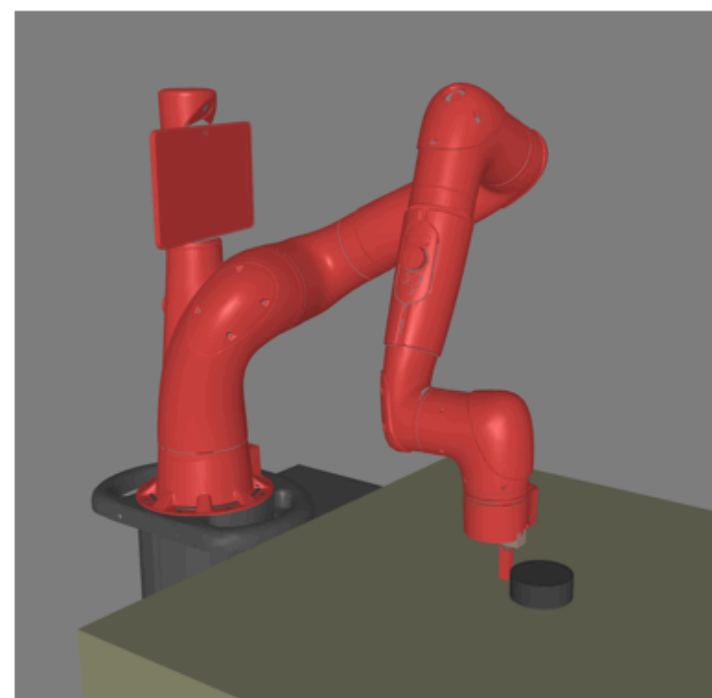
# EXPERIMENTS



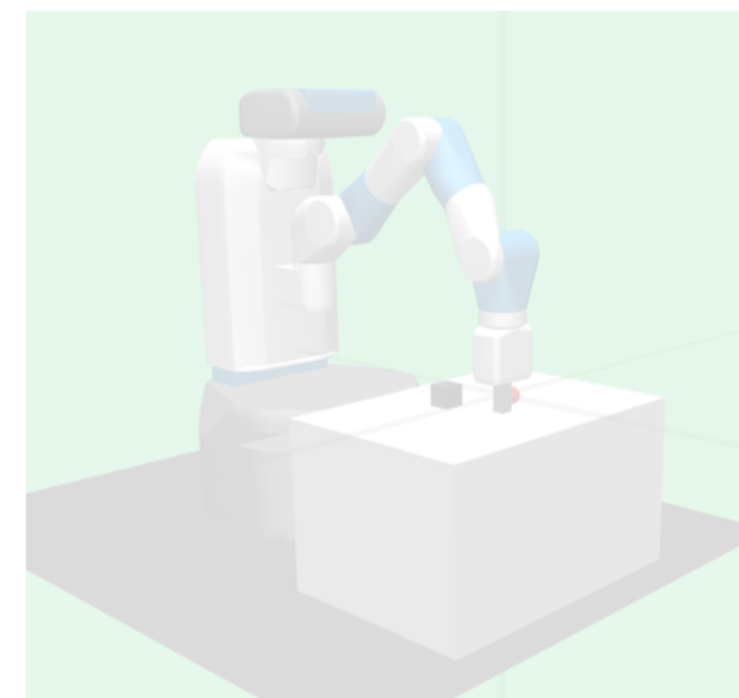
**(a) Box 2D**



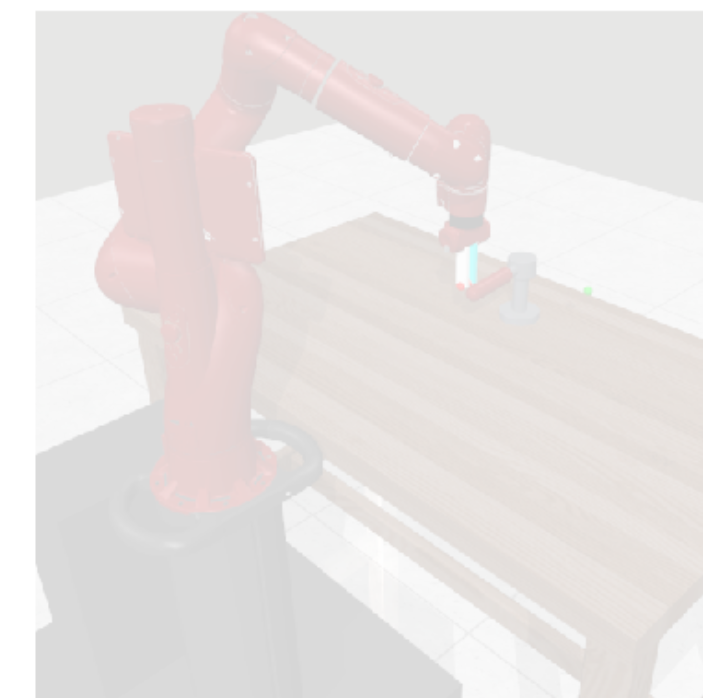
**(b) Ant**



**(c) Sawyer Push**



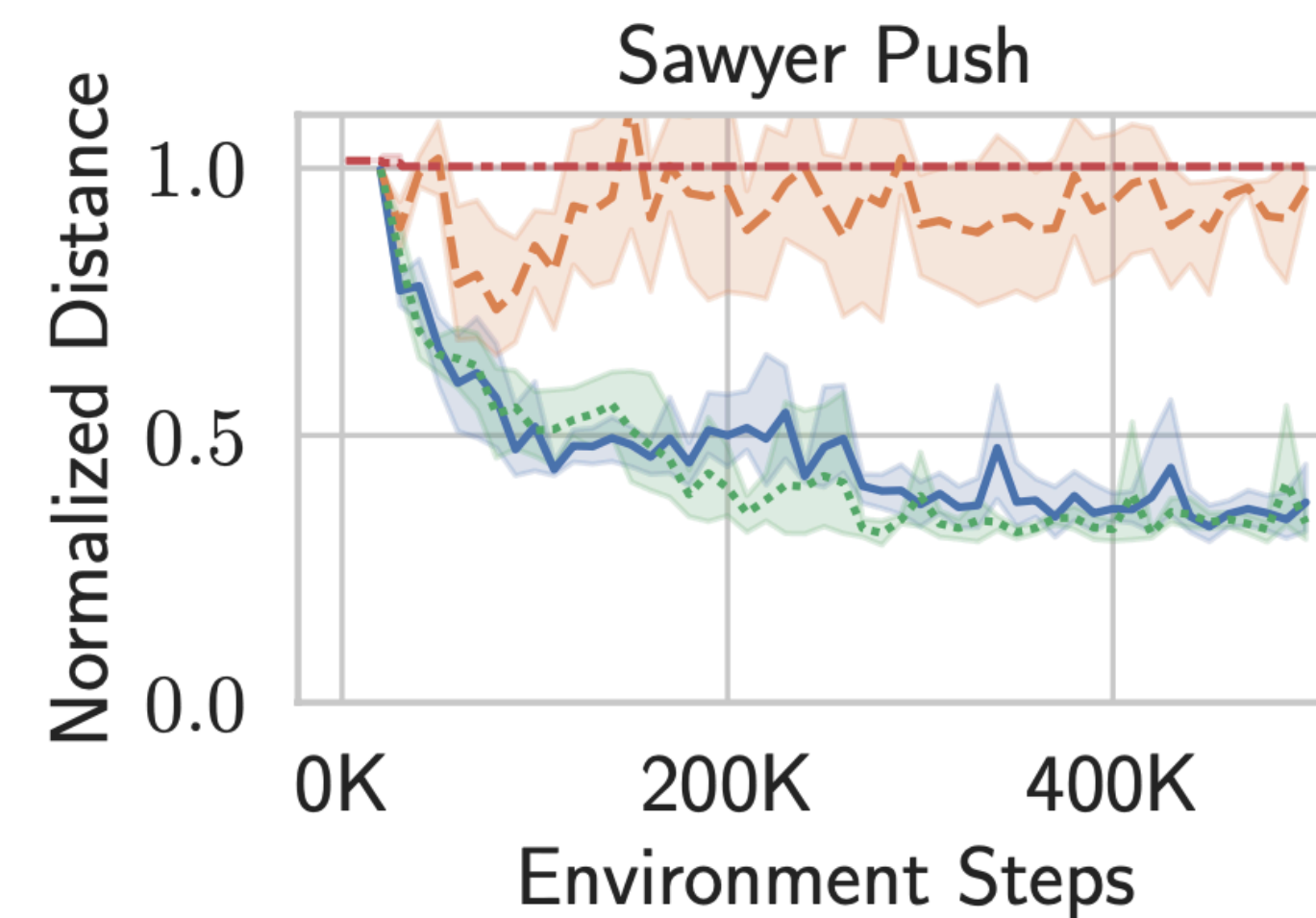
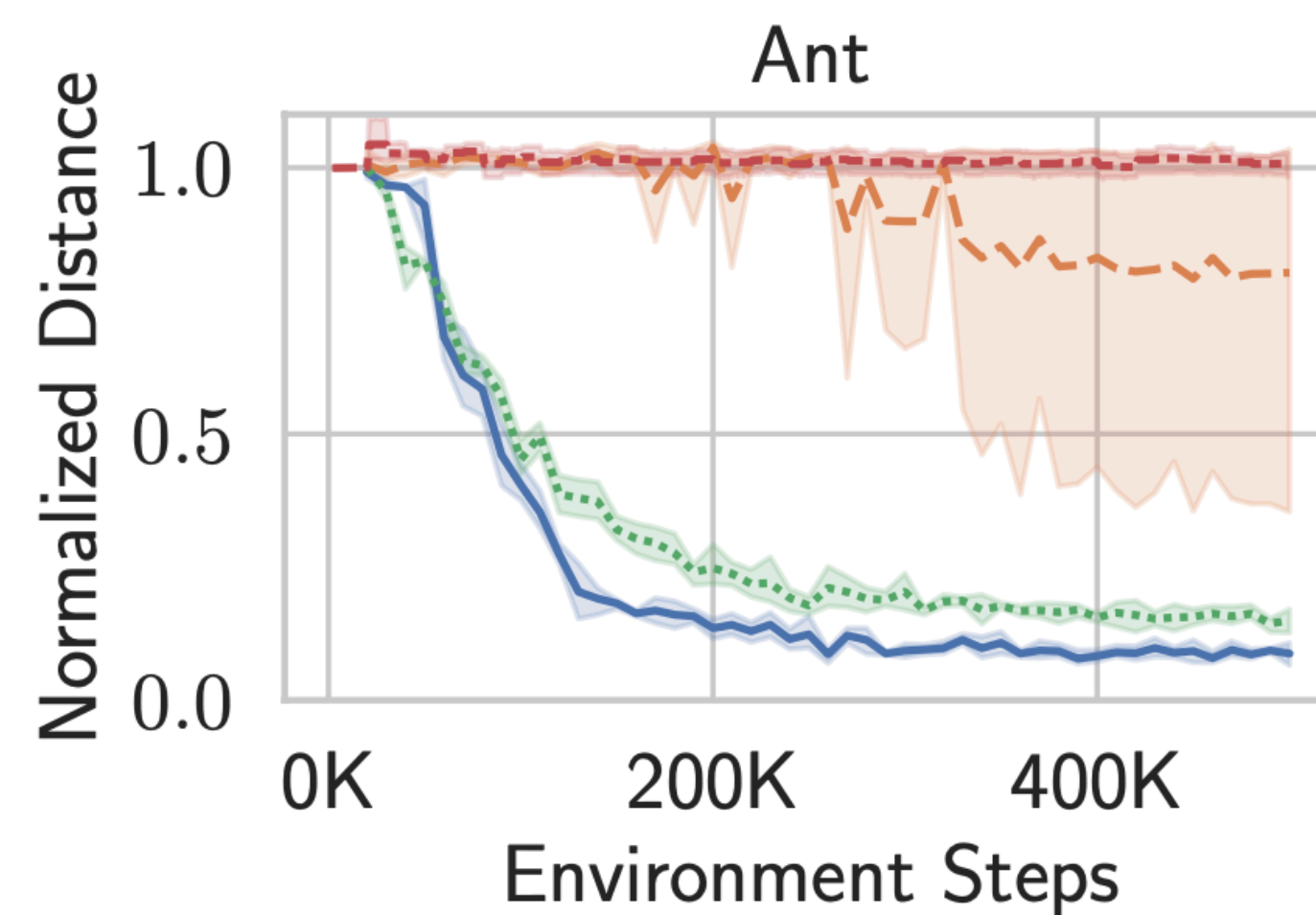
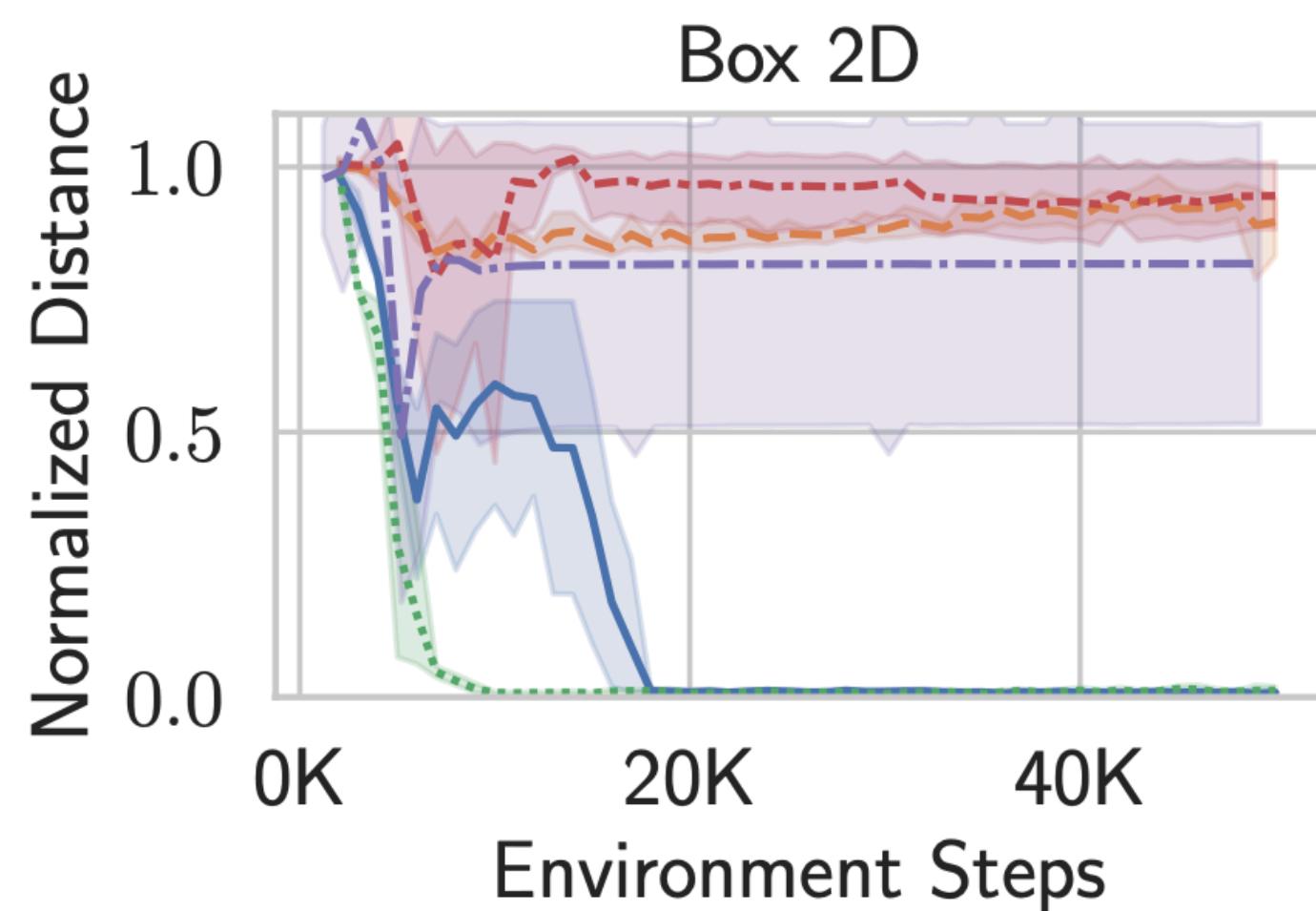
**(d) Fetch Push**



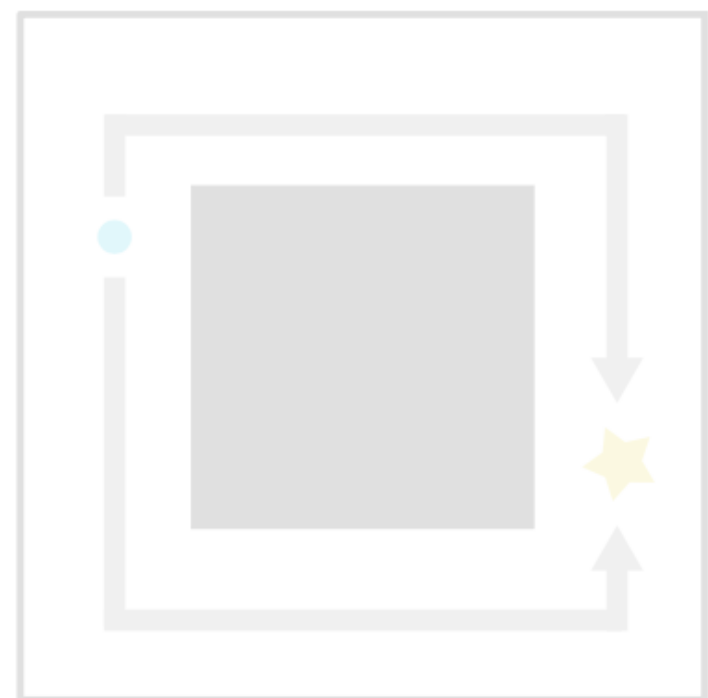
**(e) Faucet**



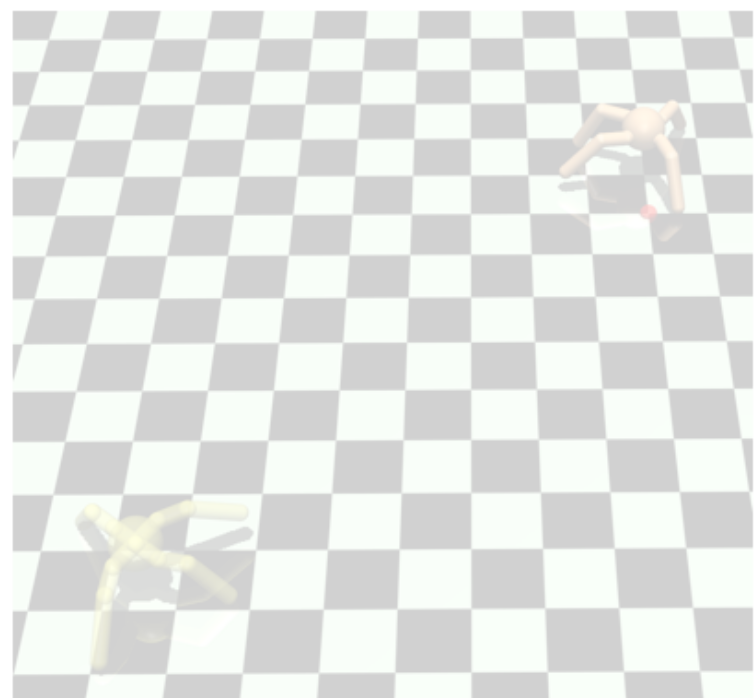
**(f) Window**



# EXPERIMENTS



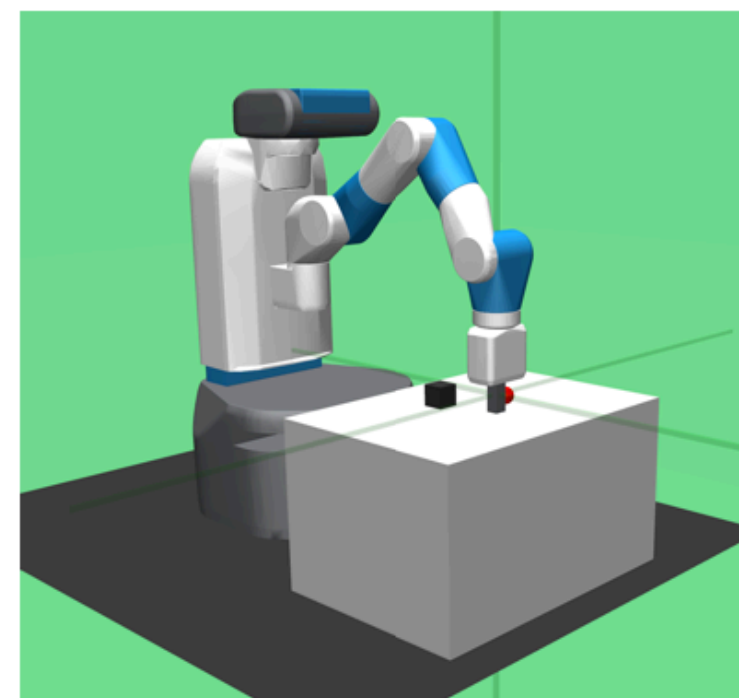
(a) Box 2D



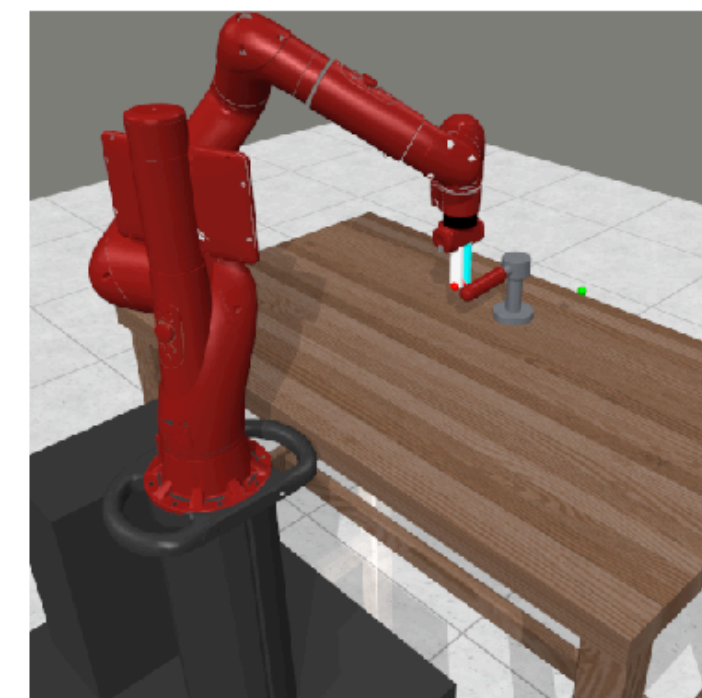
(b) Ant



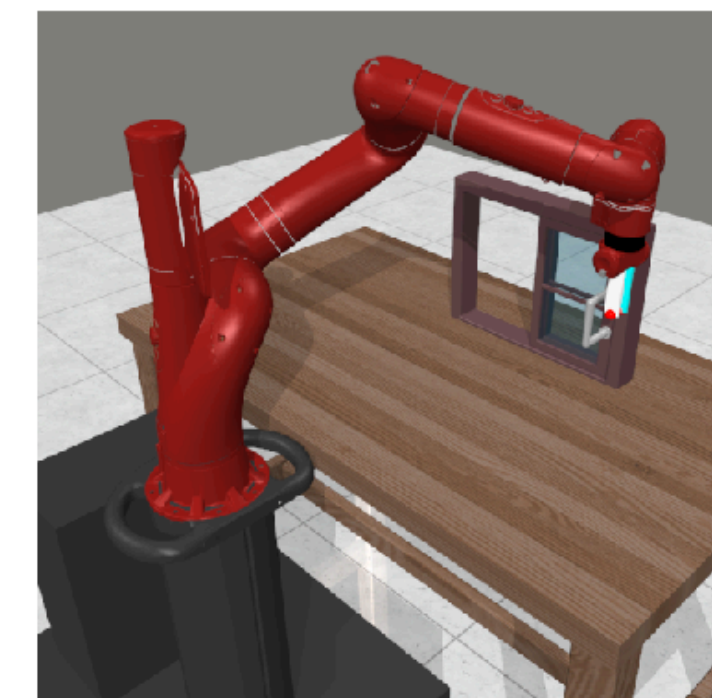
(c) Sawyer Push



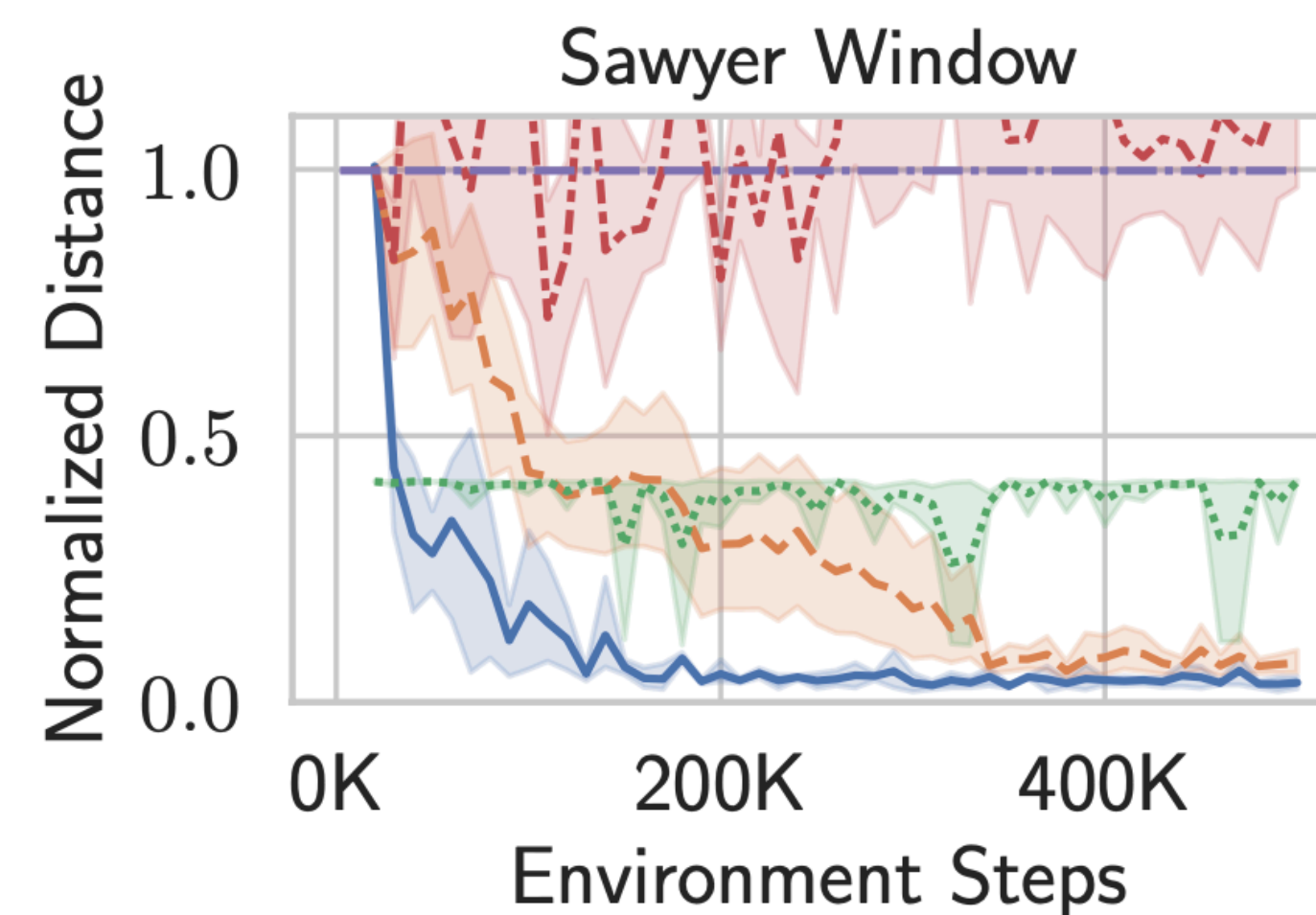
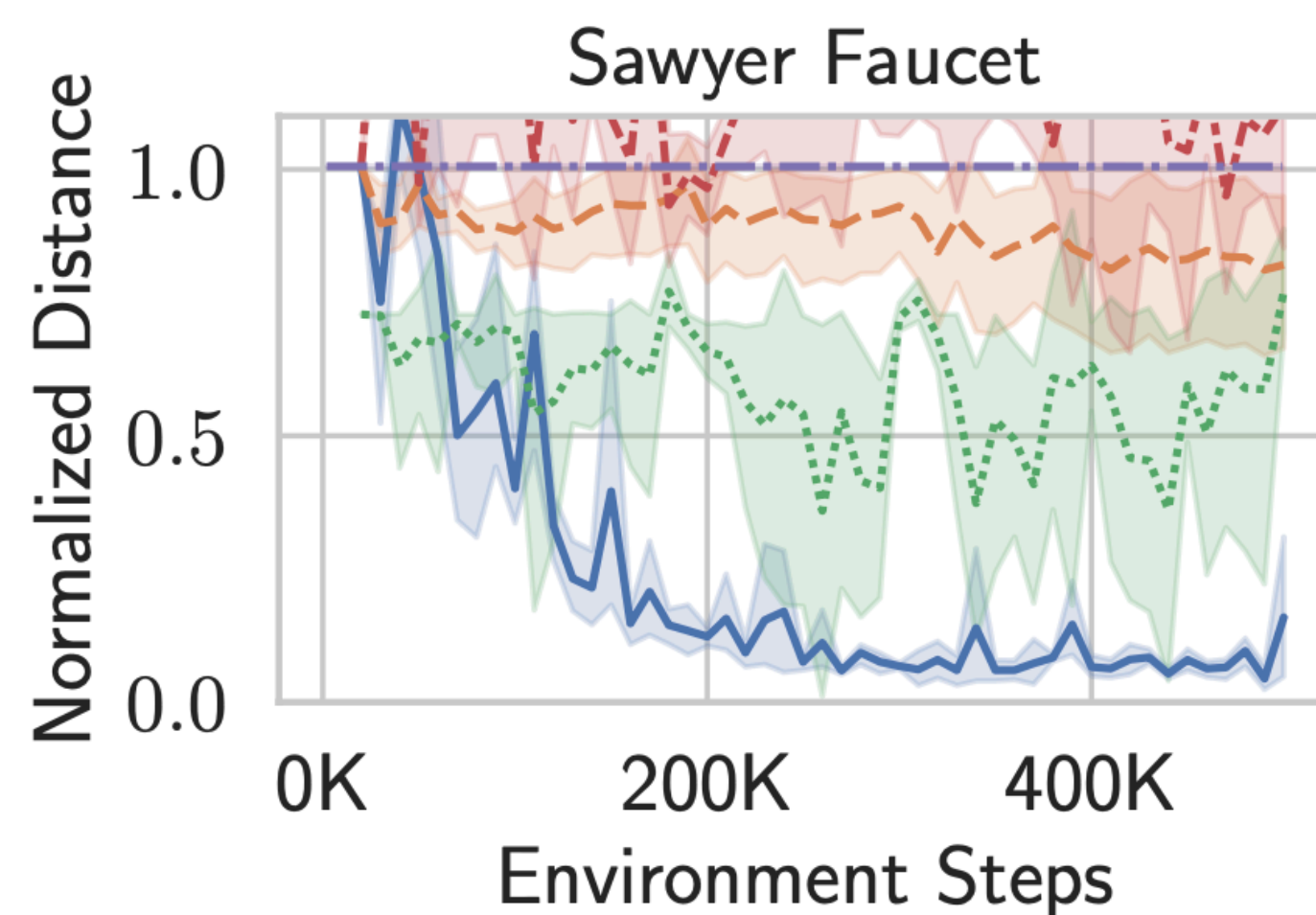
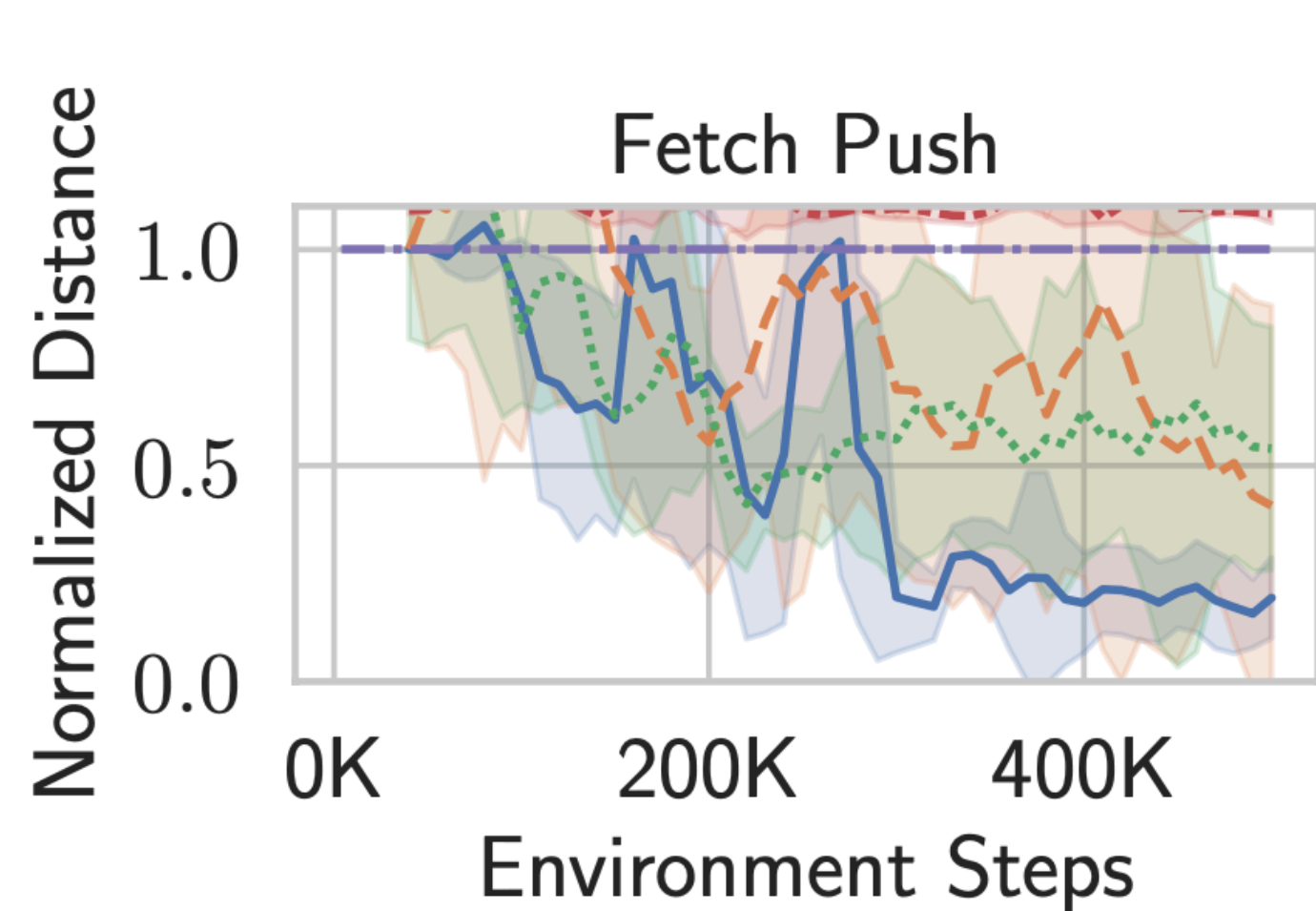
(d) Fetch Push



(e) Faucet



(f) Window



## Effect of **dynamic discount factor** and **learned dynamics**

Env	ODAC	fixed $\hat{p}_d$	fixed $q_T$	fixed $q_T, \hat{p}_d$
2D	1.7 (1.20)	1.2 (0.14)	1.0 (0.24)	1.3 (0.29)
Ant	9 (0.48)	11 (0.57)	12 (0.41)	13 (0.20)
Push	35 (2.7)	34 (1.5)	37 (1.5)	38 (3.1)
Fetch	19 (6)	15 (3)	53 (13)	66 (15)
Window	5.4 (0.62)	5.0 (0.62)	7.9 (0.71)	6.0 (0.12)
Faucet	13 (4.2)	15 (3.3)	37 (8.3)	38 (7.2)



## Effect of **dynamic discount factor** and **learned dynamics**

- Fixing the dynamics model works well

Env	ODAC	fixed $\hat{p}_d$	fixed $q_T$	fixed $q_T, \hat{p}_d$
2D	1.7 (1.20)	1.2 (0.14)	1.0 (0.24)	1.3 (0.29)
Ant	9 (0.48)	11 (0.57)	12 (0.41)	13 (0.20)
Push	35 (2.7)	34 (1.5)	37 (1.5)	38 (3.1)
Fetch	19 (6)	15 (3)	53 (13)	66 (15)
Window	5.4 (0.62)	5.0 (0.62)	7.9 (0.71)	6.0 (0.12)
Faucet	13 (4.2)	15 (3.3)	37 (8.3)	38 (7.2)

## Effect of **dynamic discount factor** and **learned dynamics**

- Fixing the dynamics model works well
- Fixing the discount factor deteriorates performance

Env	ODAC	fixed $\hat{p}_d$	fixed $q_T$	fixed $q_T, \hat{p}_d$
2D	1.7 (1.20)	1.2 (0.14)	1.0 (0.24)	1.3 (0.29)
Ant	9 (0.48)	11 (0.57)	12 (0.41)	13 (0.20)
Push	35 (2.7)	34 (1.5)	37 (1.5)	38 (3.1)
Fetch	19 (6)	15 (3)	53 (13)	66 (15)
Window	5.4 (0.62)	5.0 (0.62)	7.9 (0.71)	6.0 (0.12)
Faucet	13 (4.2)	15 (3.3)	37 (8.3)	38 (7.2)

# MAIN TAKEAWAYS

1. RL can be derived from **probabilistic inference** **without** access to a pre-specified **reward function**.

# MAIN TAKEAWAYS

1. RL can be derived from **probabilistic inference without** access to a pre-specified **reward function**.
2. **Outcome-driven** variational inference can be formulated as a **temporal-difference algorithm**,

# MAIN TAKEAWAYS

1. RL can be derived from **probabilistic inference without** access to a pre-specified **reward function**.
2. **Outcome-driven** variational inference can be formulated as a **temporal-difference algorithm**,
3. and yields a **dense reward function** that can be estimated from environment interactions.

# THANK YOU!



**TIM G. J. RUDNER\***



**VITCHYR H. PONG\***



ROWAN McALLISTER



YARIN GAL



SERGEY LEVINE

CORRESPONDENCE: `tim.rudner@cs.ox.ac.uk` & `vitchyr@berkeley.edu`

PROJECT WEBSITE: `https://sites.google.com/view/od-ac`