

# Model-based Reinforcement Learning

Aidan Scannell

Finnish Center for Artificial Intelligence (FCAI)

Aalto University

17th July 2024

FCAI

Slides available here

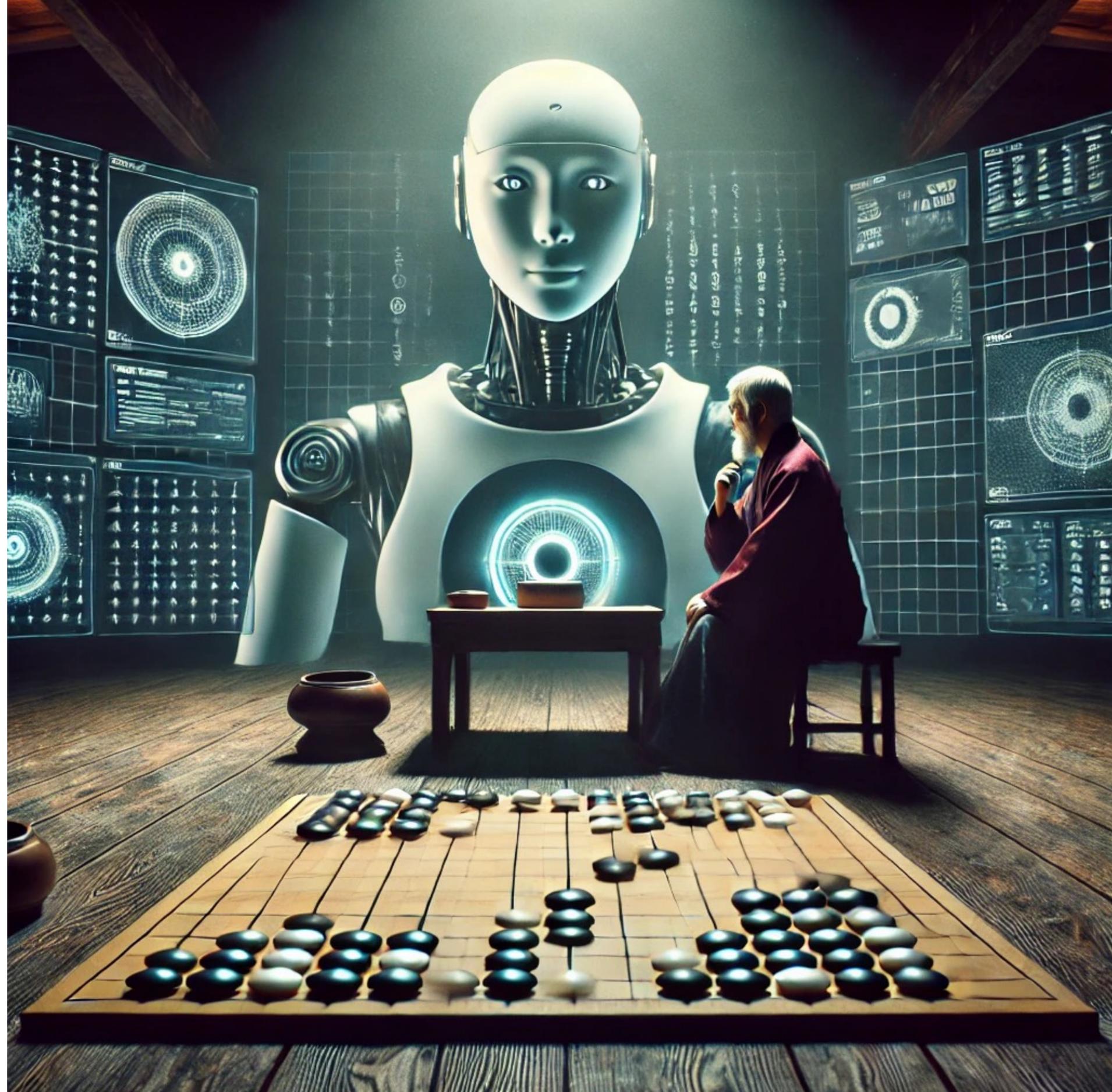


# AlphaGo

Model-based reasoning for games

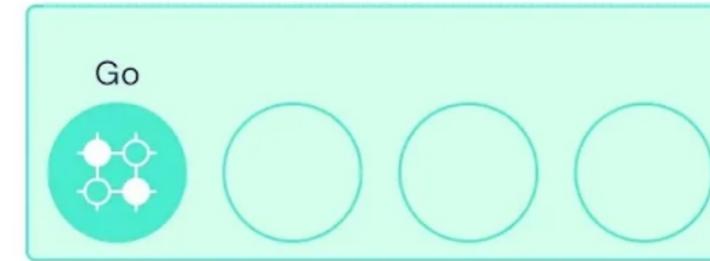
Silver et al. (2016). Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587), 484.

FCAI

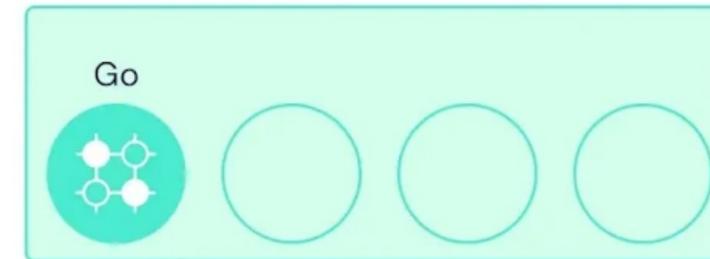


# AlphaGo

## Model-based reasoning for games



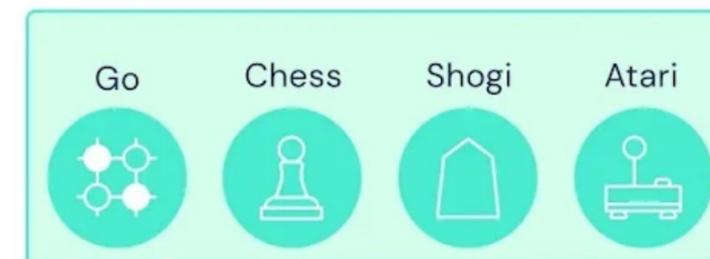
**AlphaGo** becomes the first program to master Go using neural networks and tree search (Jan 2016, Nature)



**AlphaGo Zero** learns to play completely on its own, without human knowledge (Oct 2017, Nature)



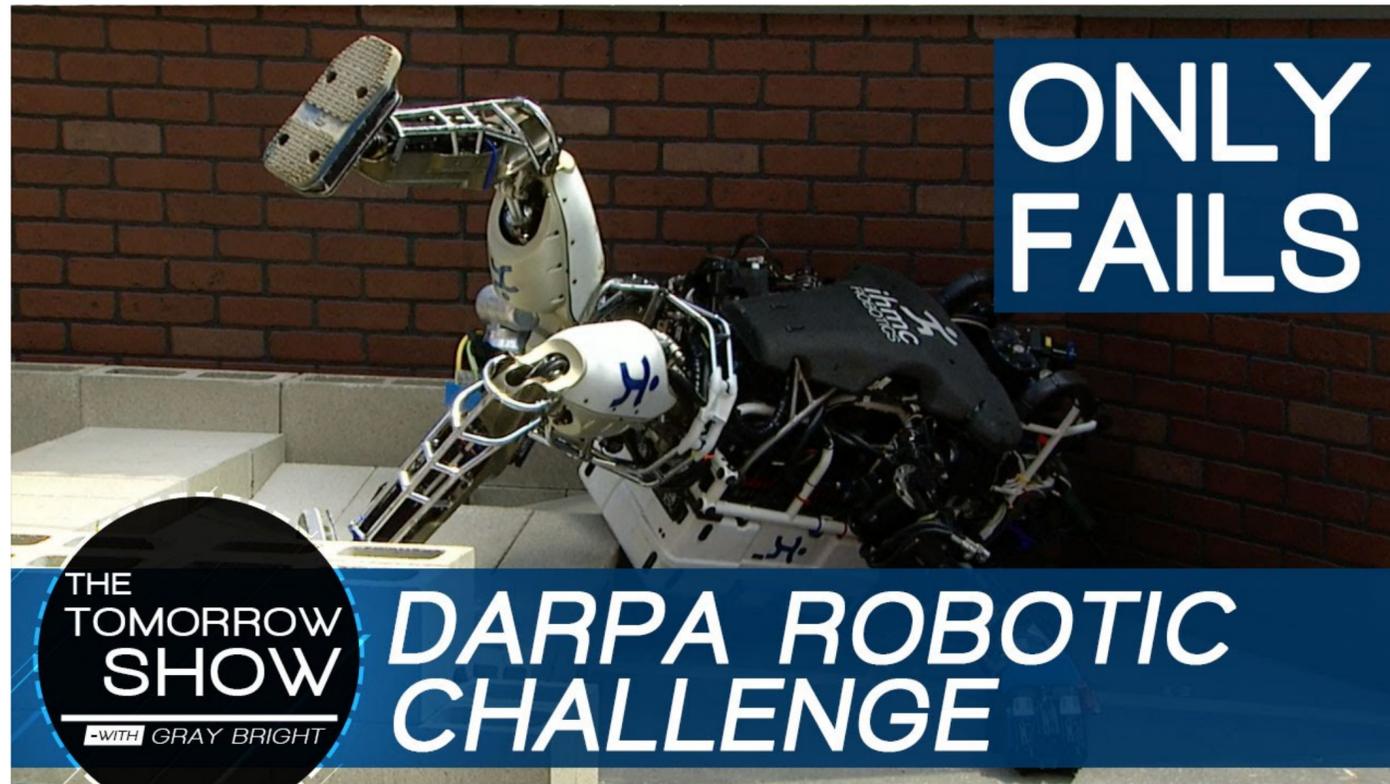
**AlphaZero** masters three perfect information games using a single algorithm for all games (Dec 2018, Science)



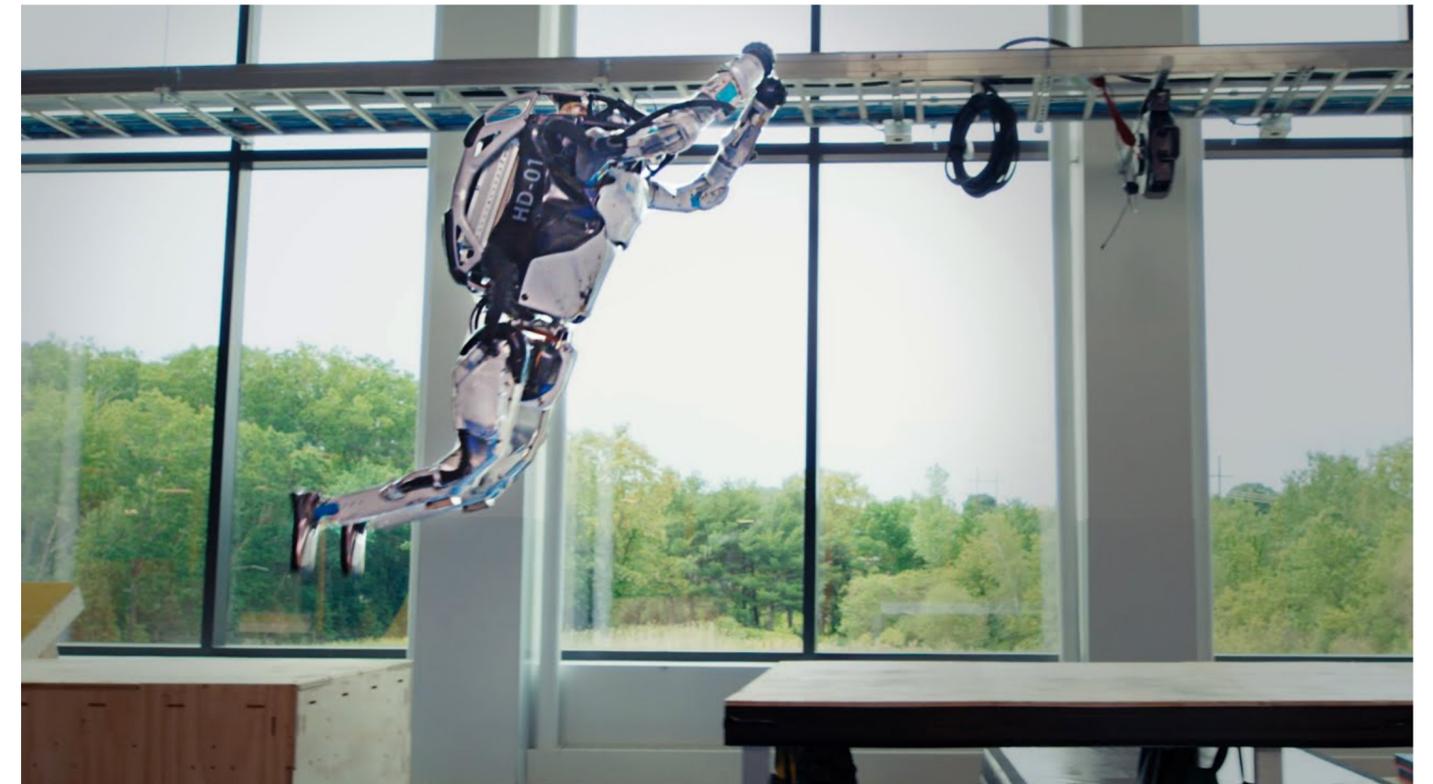
**MuZero** learns the rules of the game, allowing it to also master environments with unknown dynamics. (Dec 2020, Nature)

Silver et al. (2016). Mastering the game of Go with deep neural networks and tree search. Nature, 529(7587), 484.

# Machine Learning for Robotics

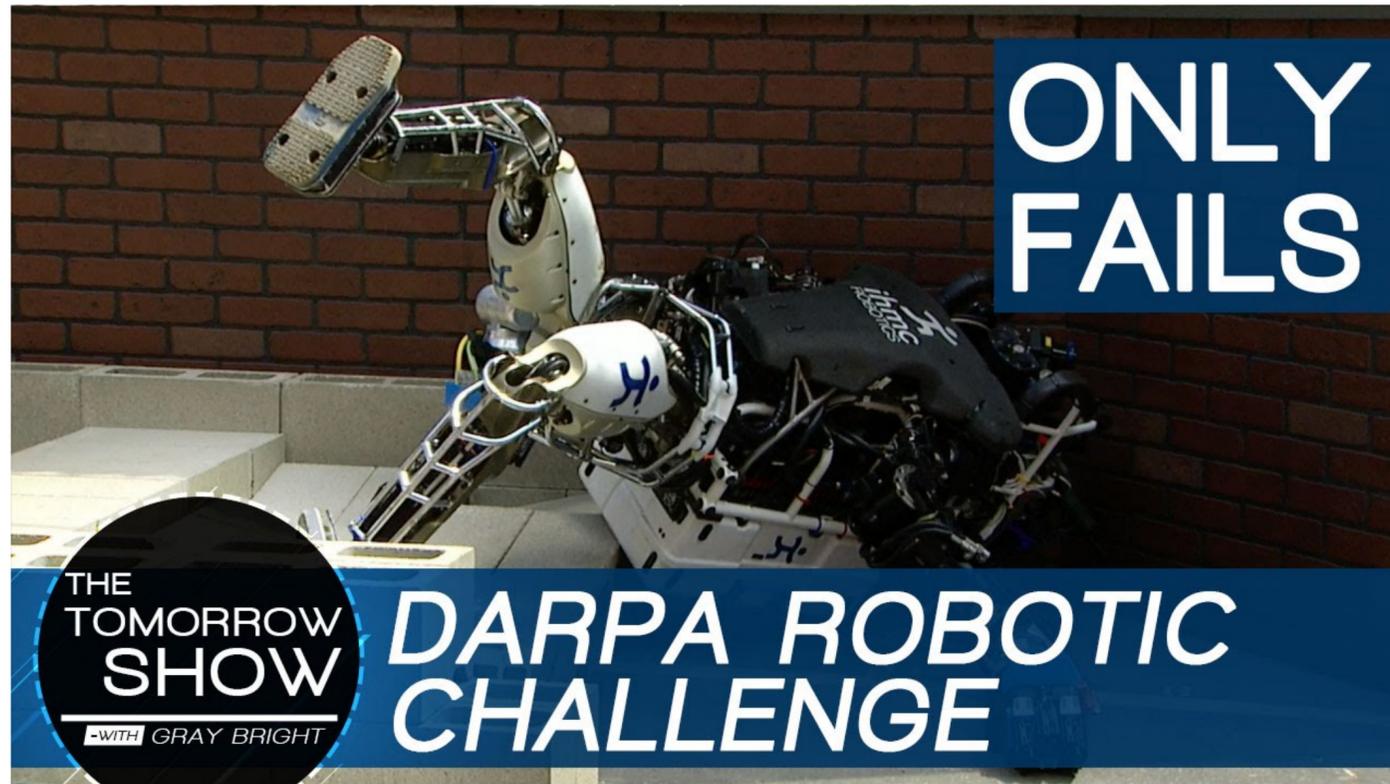


DARPA Robotics Challenge 2015

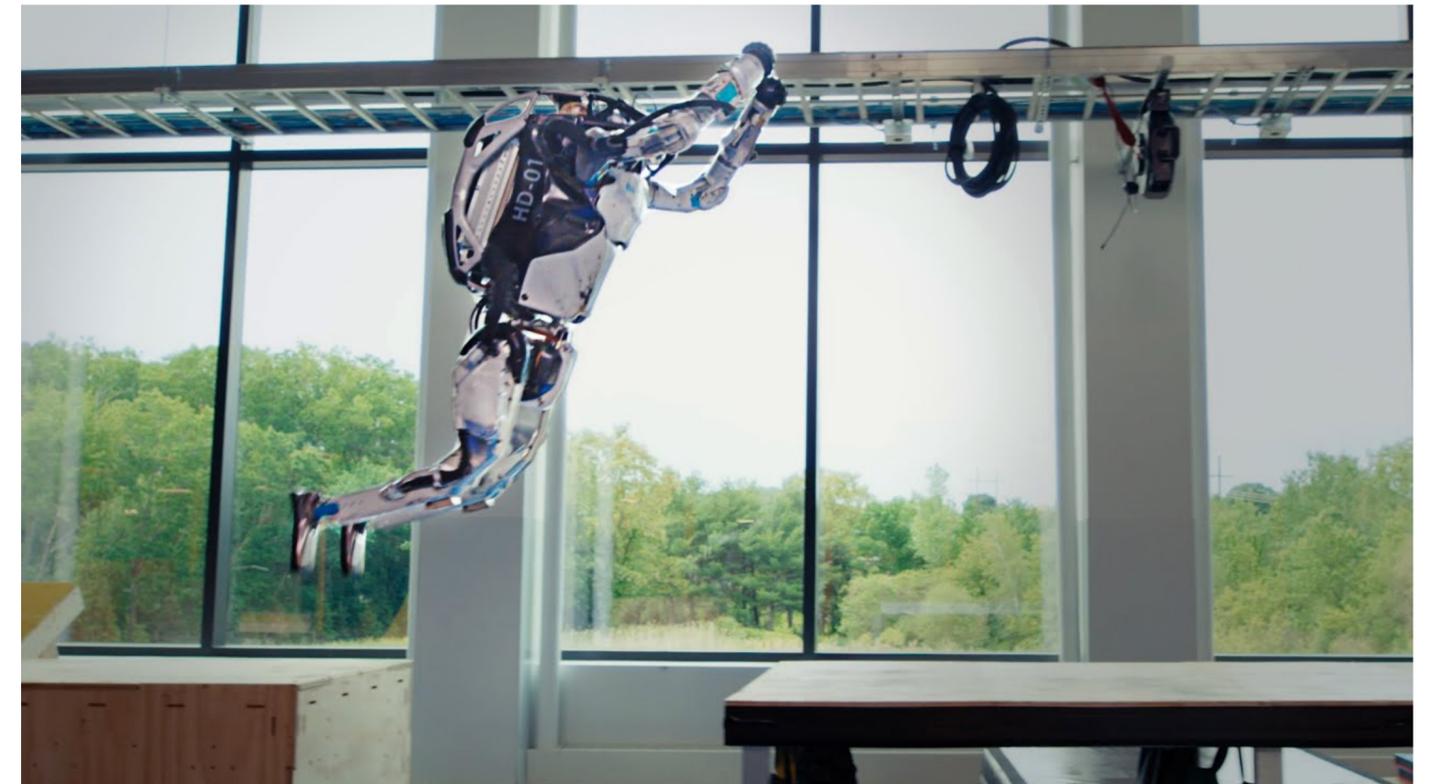


Boston Dynamics Atlas - Partners in Parkour

# Machine Learning for Robotics

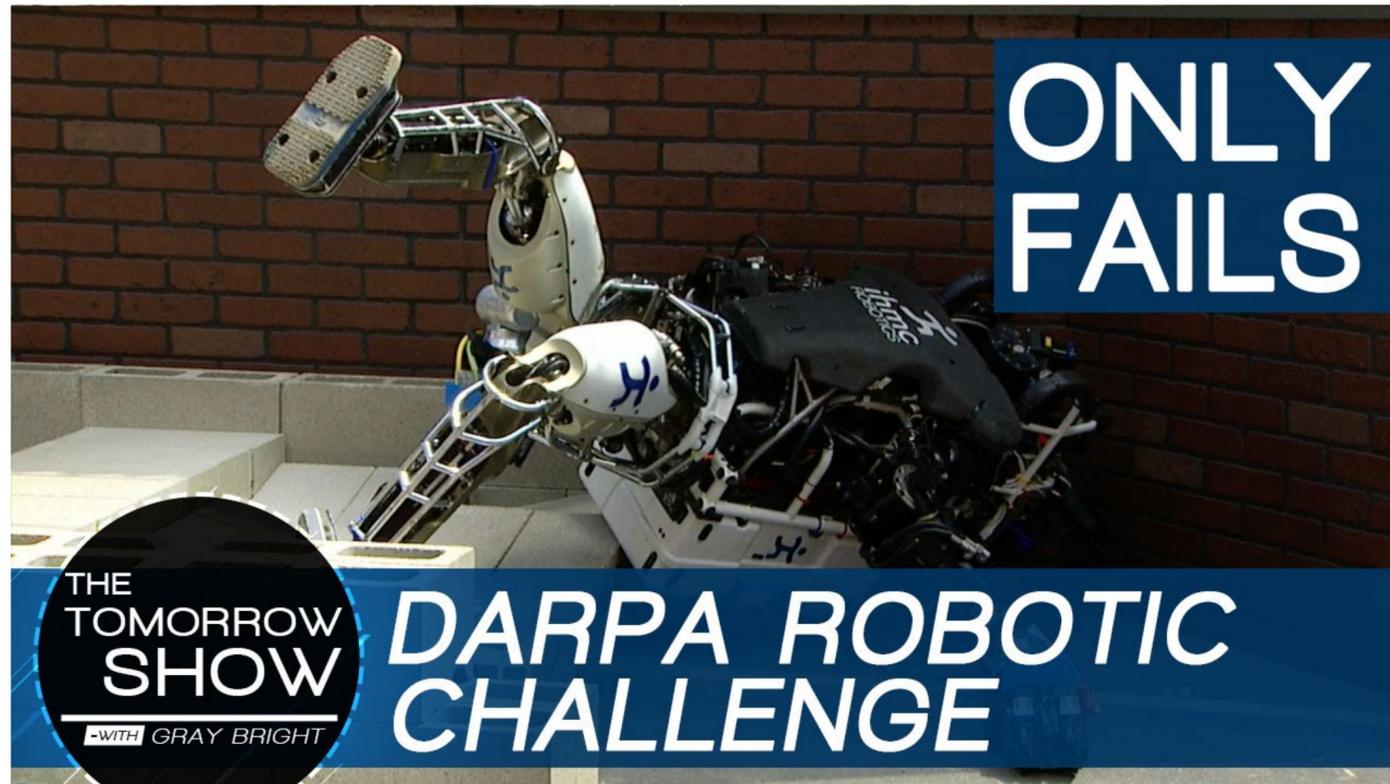


DARPA Robotics Challenge 2015

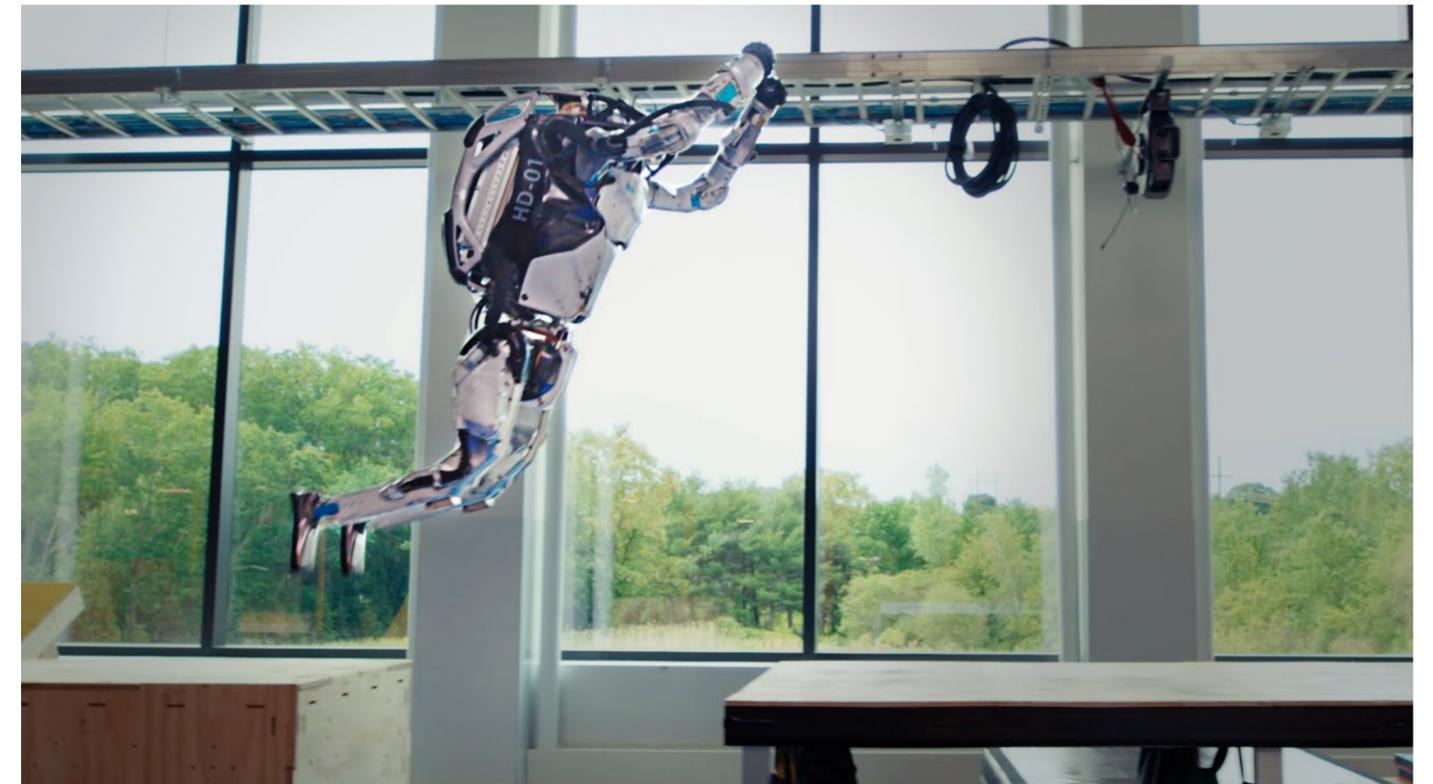


Boston Dynamics Atlas - Partners in Parkour

# Machine Learning for Robotics



DARPA Robotics Challenge 2015



Boston Dynamics Atlas - Partners in Parkour

# Learning Objectives

# Learning Objectives

Understand

# Learning Objectives

Understand

1. What a “model” is in model-based RL

# Learning Objectives

Understand

1. What a “model” is in model-based RL
2. How a “model” can aid decision making

# Learning Objectives

Understand

1. What a “model” is in model-based RL
2. How a “model” can aid decision making
3. Differences between background and decision-time planning

# Learning Objectives

Understand

1. What a “model” is in model-based RL
2. How a “model” can aid decision making
3. Differences between background and decision-time planning
4. Decision-time planning strategies for continuous actions

# Learning Objectives

## Understand

1. What a “model” is in model-based RL
2. How a “model” can aid decision making
3. Differences between background and decision-time planning
4. Decision-time planning strategies for continuous actions
5. Sources of uncertainty in model-based RL

# Learning Objectives

## Understand

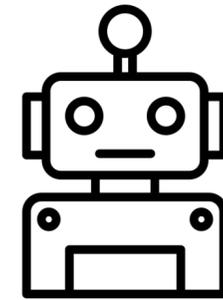
1. What a “model” is in model-based RL
2. How a “model” can aid decision making
3. Differences between background and decision-time planning
4. Decision-time planning strategies for continuous actions
5. Sources of uncertainty in model-based RL
6. Rationale and insights for decision-making under uncertainty

# Reinforcement Learning

## Markov Decision Process (MDP)

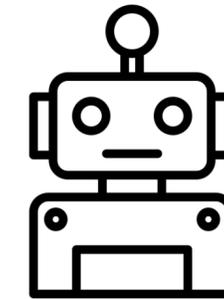
# Reinforcement Learning

## Markov Decision Process (MDP)



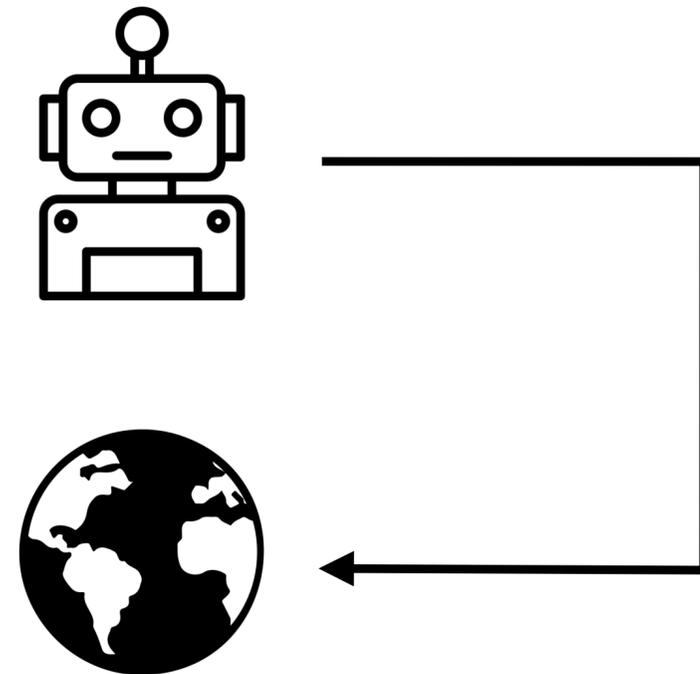
# Reinforcement Learning

## Markov Decision Process (MDP)



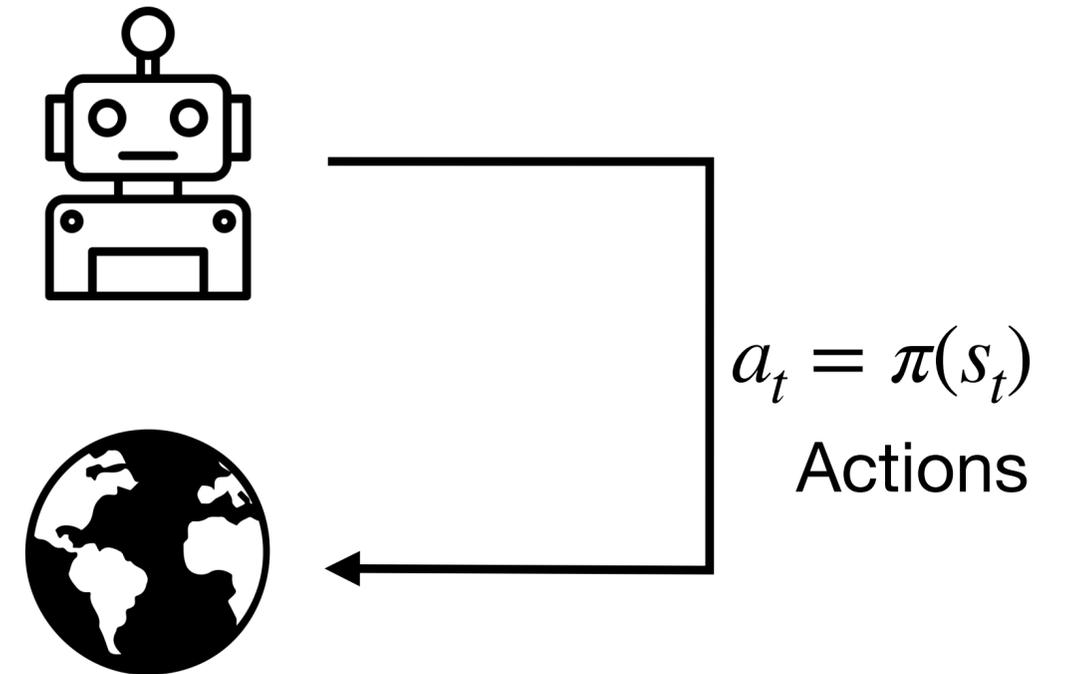
# Reinforcement Learning

## Markov Decision Process (MDP)



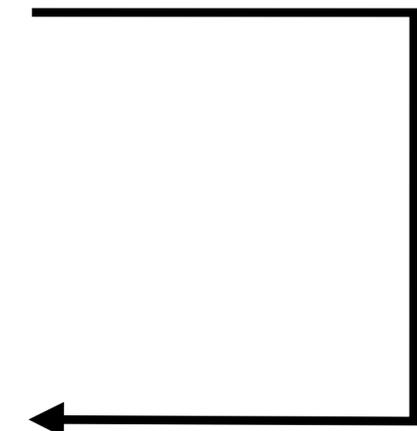
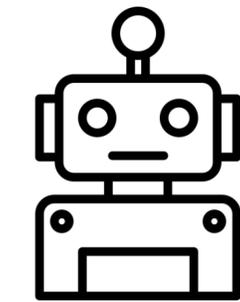
# Reinforcement Learning

## Markov Decision Process (MDP)



# Reinforcement Learning

## Markov Decision Process (MDP)



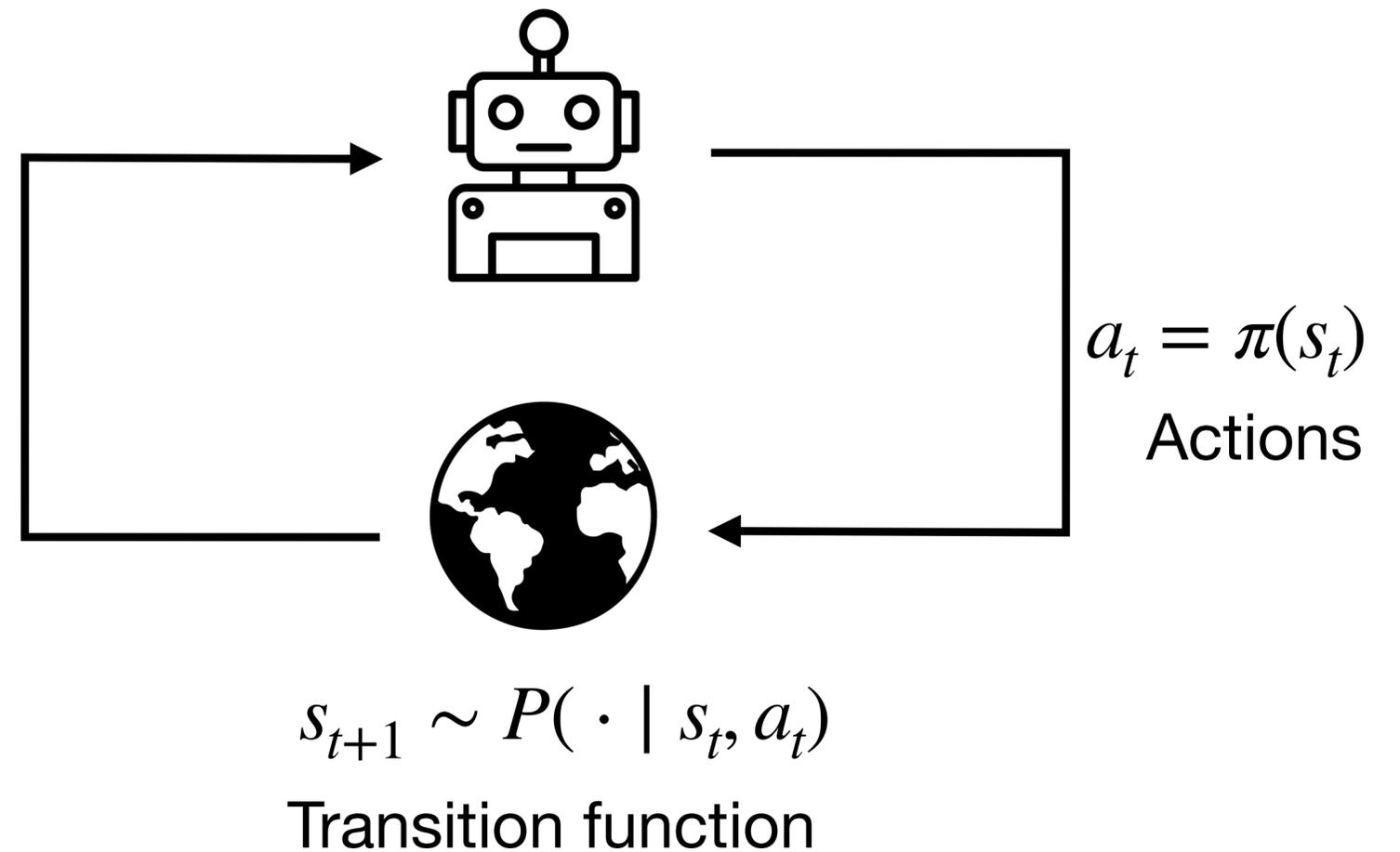
$a_t = \pi(s_t)$   
Actions

$$s_{t+1} \sim P(\cdot | s_t, a_t)$$

Transition function

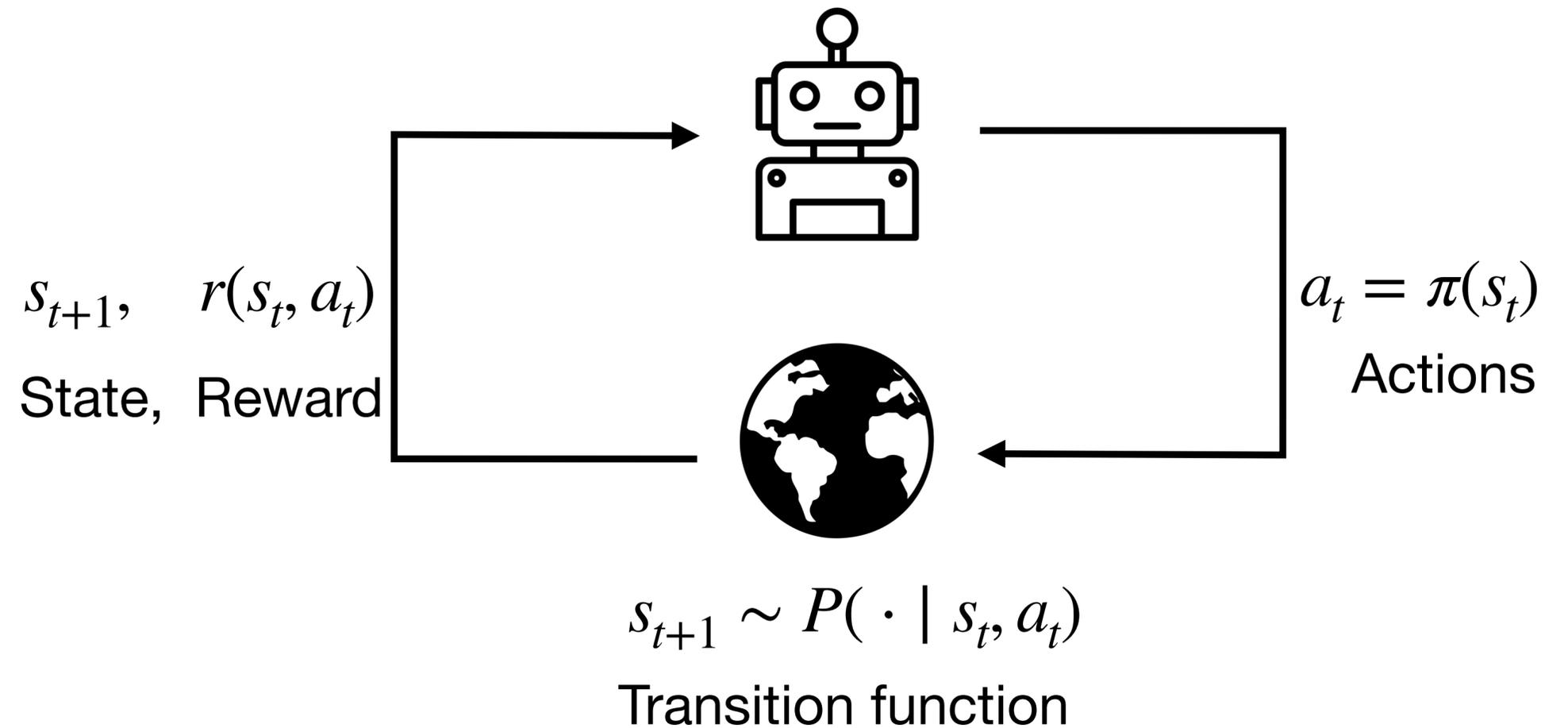
# Reinforcement Learning

## Markov Decision Process (MDP)



# Reinforcement Learning

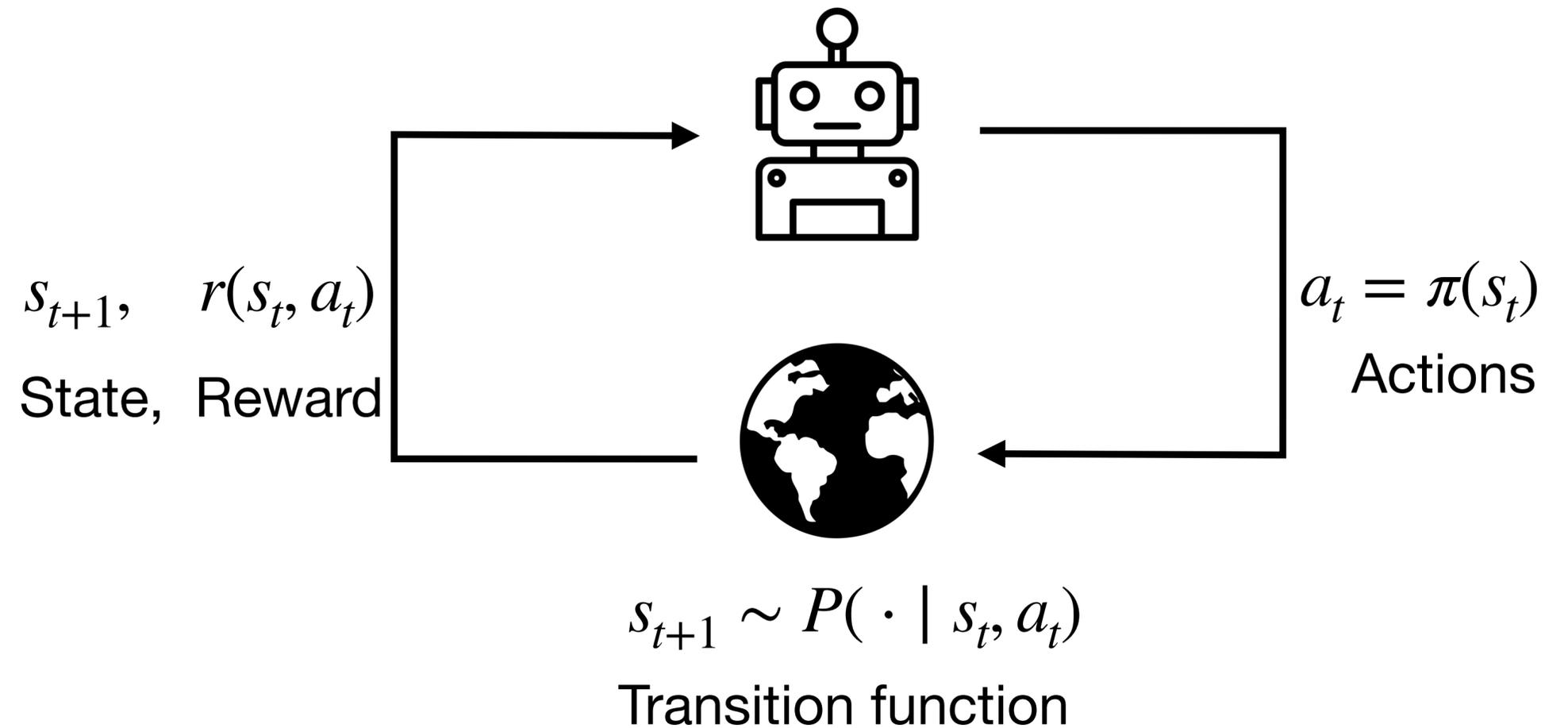
## Markov Decision Process (MDP)



# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

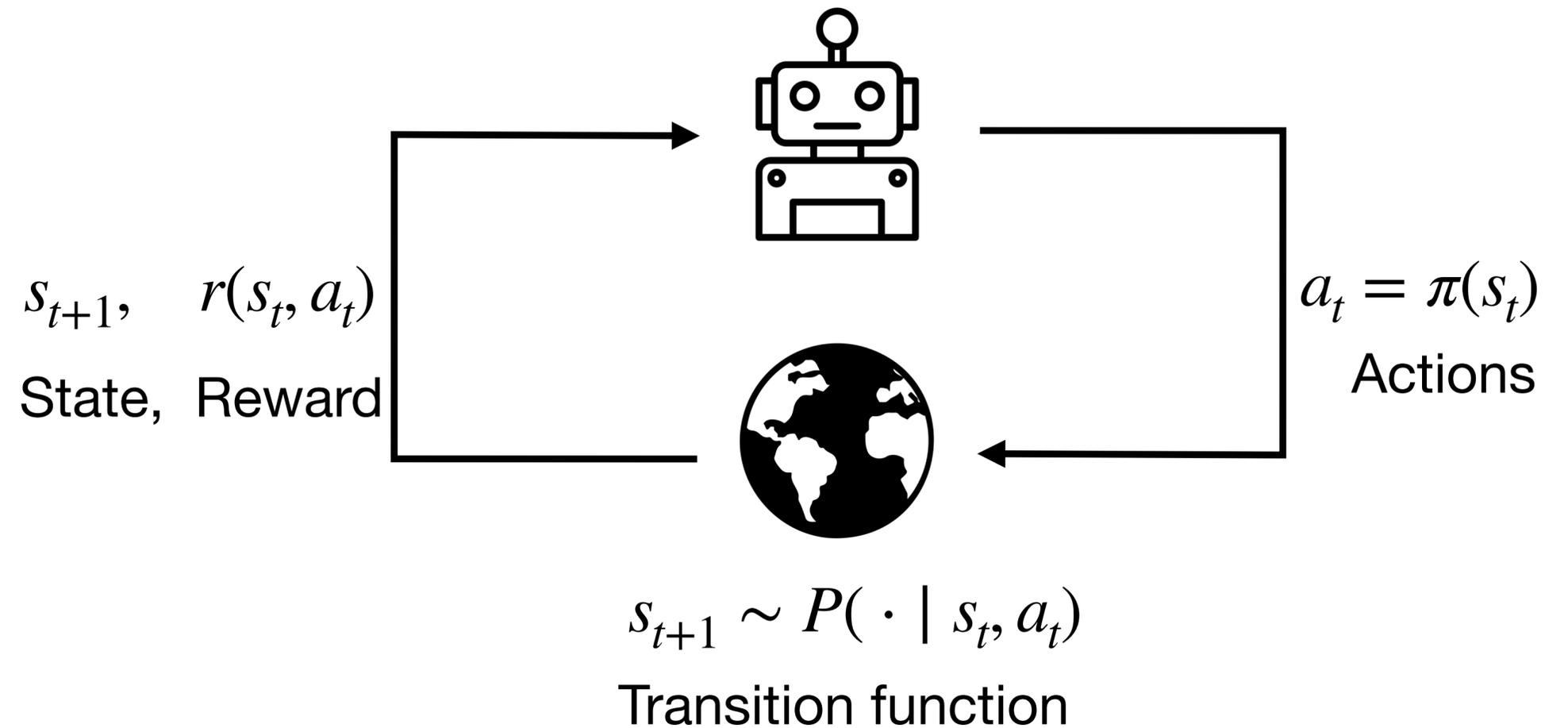


# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$



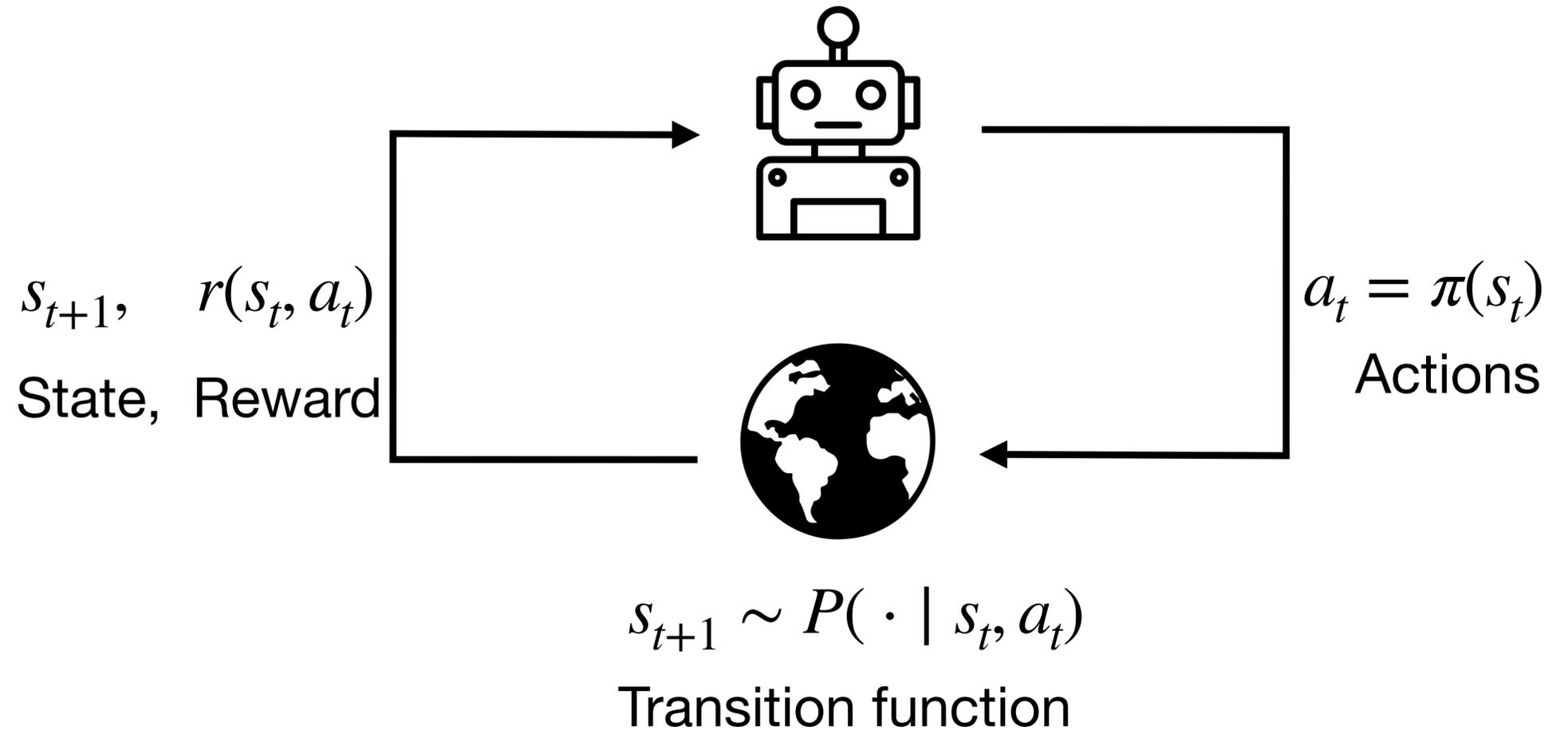
# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} | s_t, a_t)$



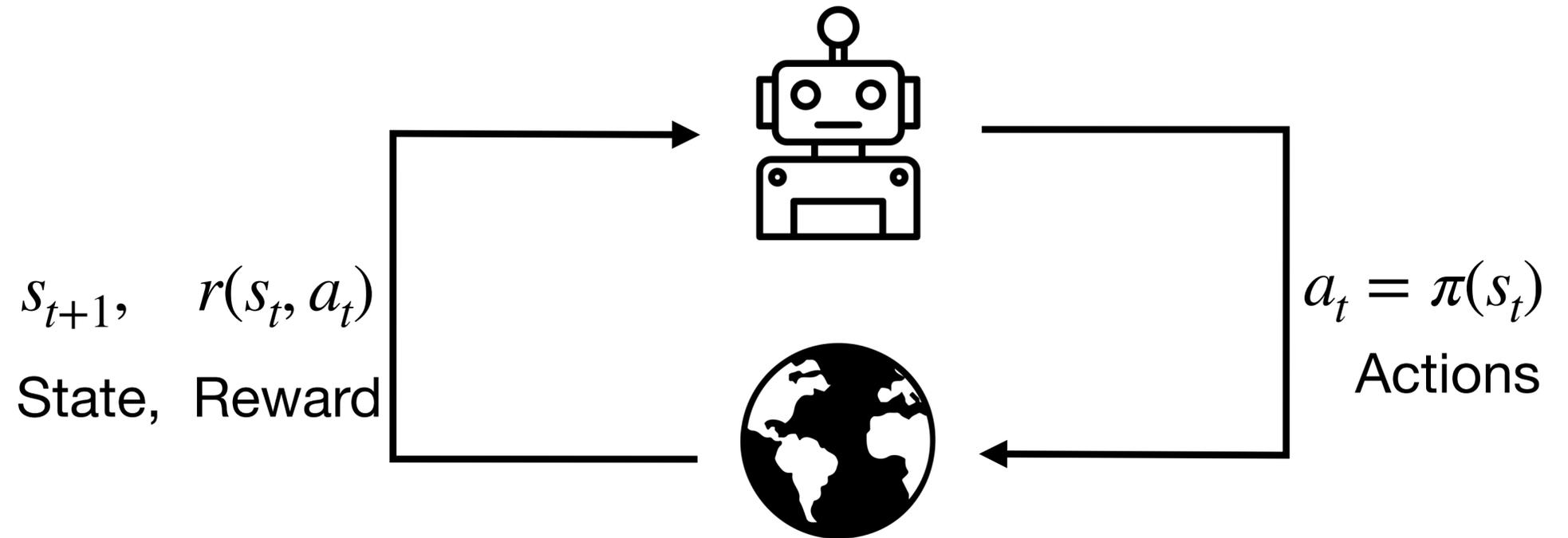
# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} | s_t, a_t)$



$$s_{t+1} \sim P(\cdot | s_t, a_t)$$

Transition function

$$s_{t+1} = f(s_t, a_t) + \epsilon_t$$

# Reinforcement Learning

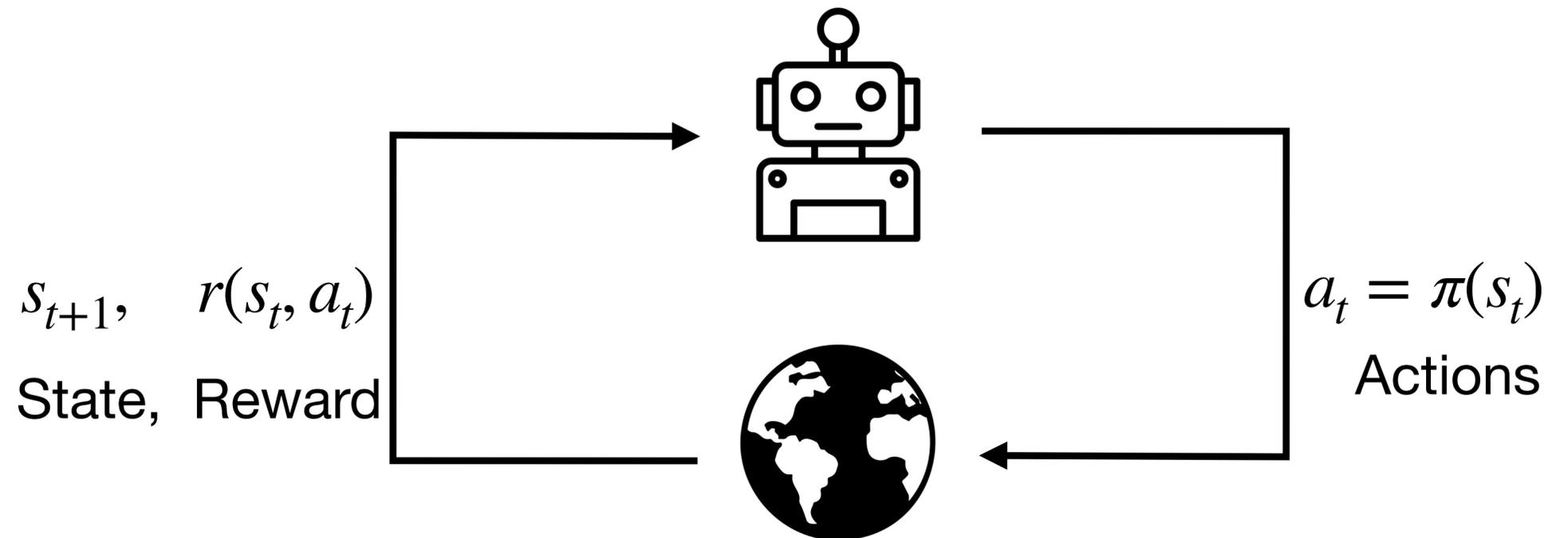
## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} | s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$



$$s_{t+1} \sim P(\cdot | s_t, a_t)$$

Transition function

$$s_{t+1} = f(s_t, a_t) + \epsilon_t$$

# Reinforcement Learning

## Markov Decision Process (MDP)

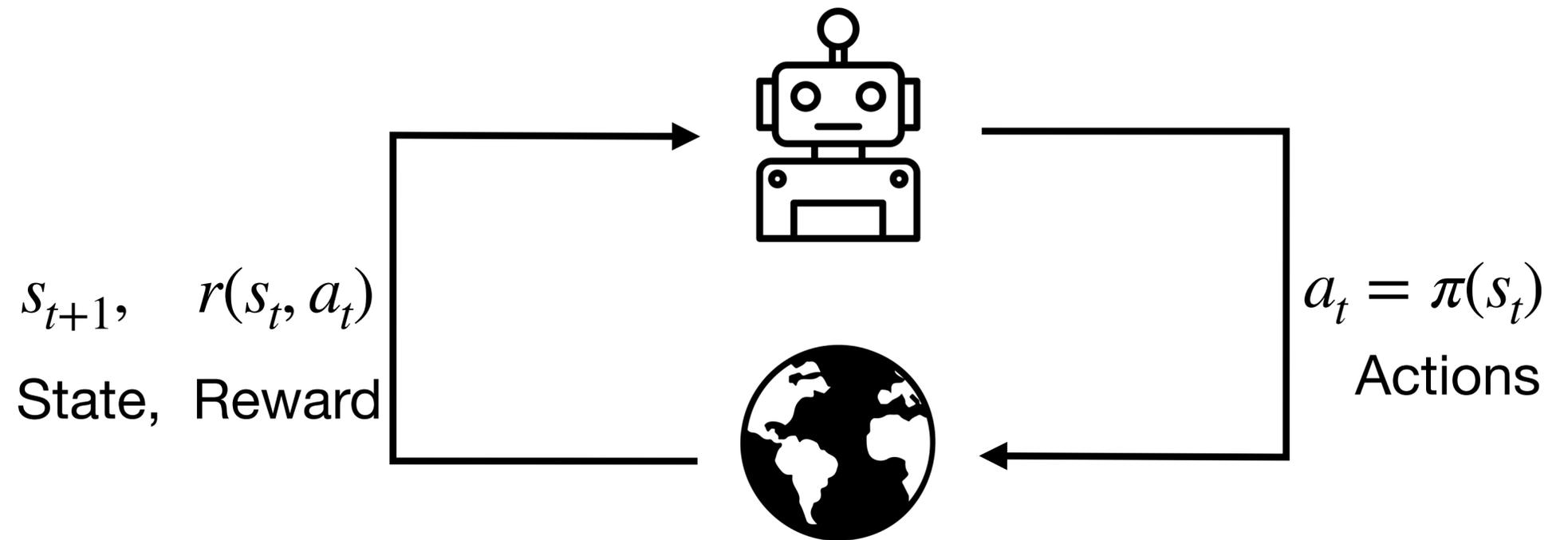
States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} | s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$



$$s_{t+1} \sim P(\cdot | s_t, a_t)$$

Transition function

$$s_{t+1} = f(s_t, a_t) + \epsilon_t$$

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

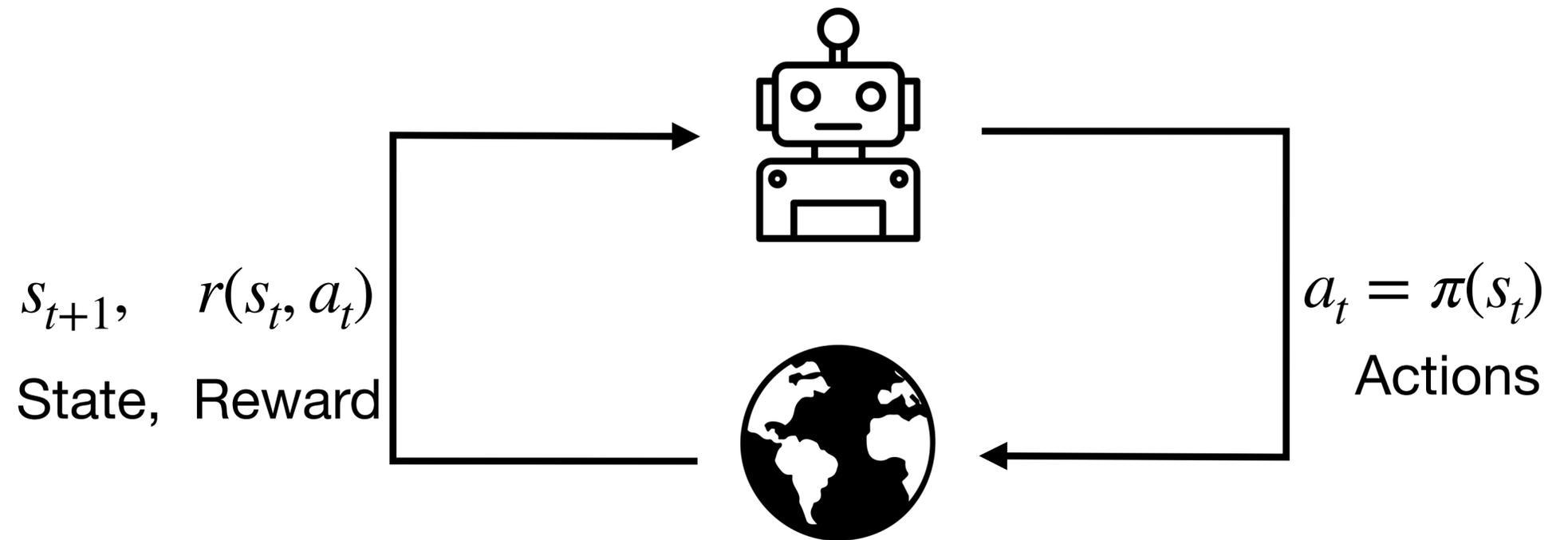
Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} | s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0,1]$



$s_{t+1}, r(s_t, a_t)$   
State, Reward

$a_t = \pi(s_t)$   
Actions

$$s_{t+1} \sim P(\cdot | s_t, a_t)$$

Transition function

$$s_{t+1} = f(s_t, a_t) + \epsilon_t$$

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

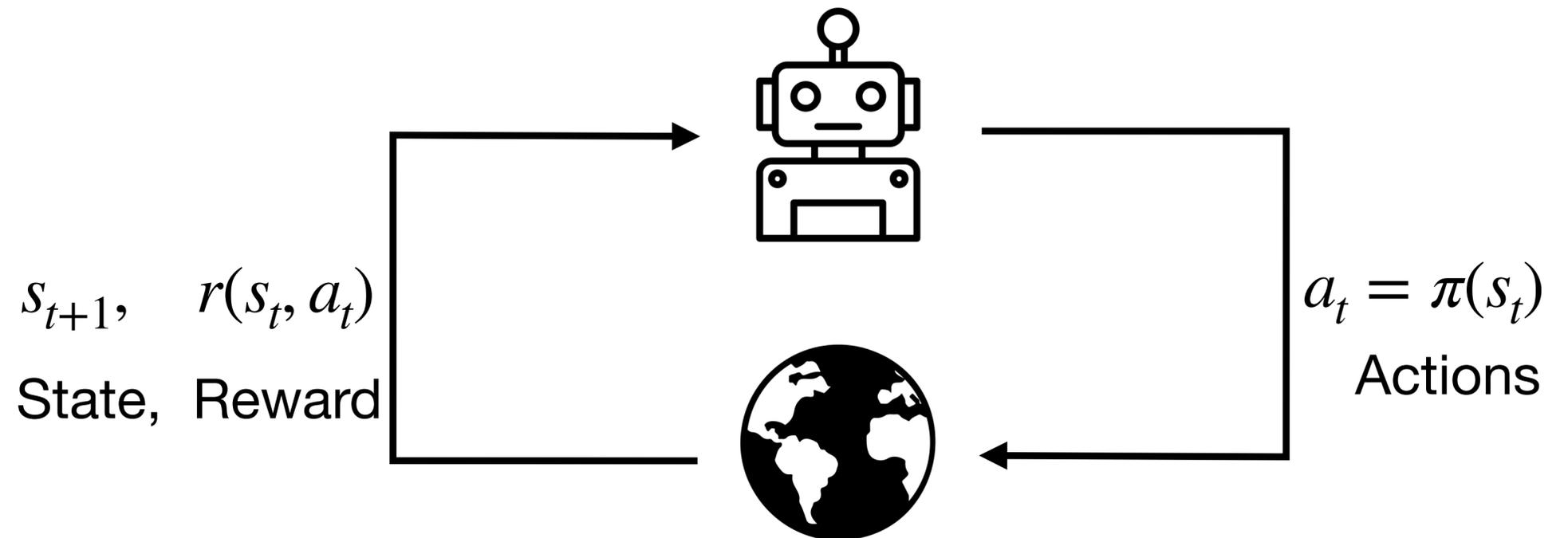
Transition function  $P(s_{t+1} | s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0,1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$



$$s_{t+1} \sim P(\cdot | s_t, a_t)$$

Transition function

$$s_{t+1} = f(s_t, a_t) + \epsilon_t$$

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} \mid s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0,1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

**Goal:**

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} \mid s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0,1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} \mid s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0, 1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

**Goal:**

$$\max_{\pi} \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} \mid s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0, 1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

**Goal:**

$$\max_{\pi} \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} \mid s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0, 1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

**Goal:**

$$\max_{\pi} \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

**Value function:**

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} \mid s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0, 1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

**Goal:**

$$\max_{\pi} \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

**Value function:**

$$V_{\pi}(s) = \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} \mid s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0, 1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

**Goal:**

$$\max_{\pi} \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

**Value function:**

$$V_{\pi}(s) = \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} \mid s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0, 1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

**Goal:**

$$\max_{\pi} \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

**Value function:**

$$V_{\pi}(s) = \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

**Action-value function (aka Q-function):**

# Reinforcement Learning

## Markov Decision Process (MDP)

States  $s \in \mathcal{S}$

Actions  $a \in \mathcal{A}$

Transition function  $P(s_{t+1} \mid s_t, a_t)$

Reward function  $r_t = r(s_t, a_t)$

Start state  $s_0$

Discount factor  $\gamma \in [0, 1]$

Policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$

**Goal:**

$$\max_{\pi} \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

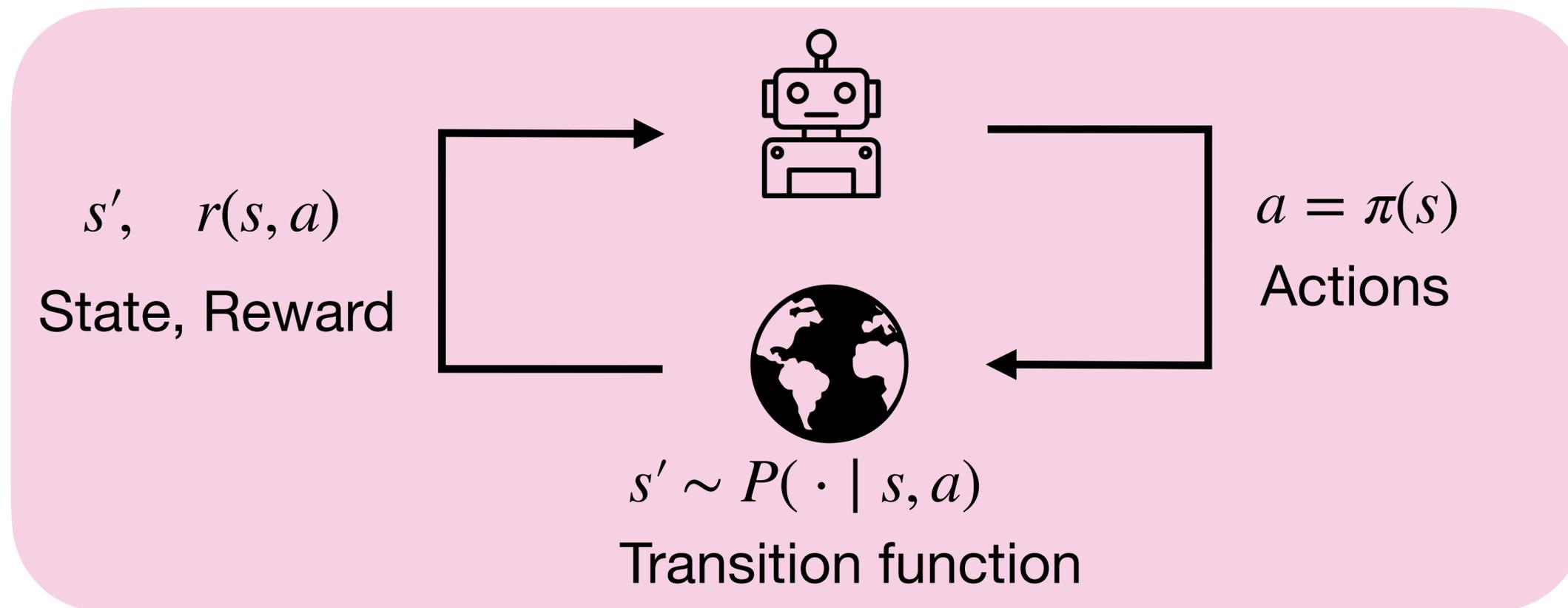
**Value function:**

$$V_{\pi}(s) = \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi \right]$$

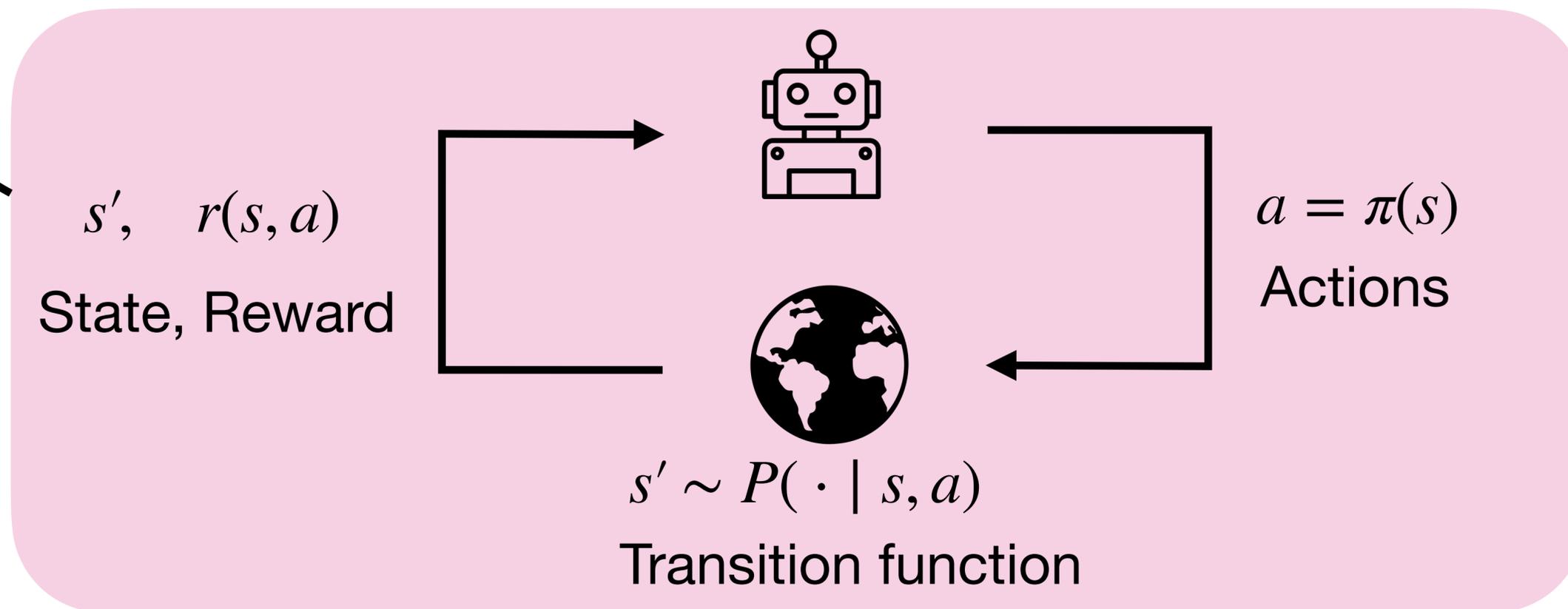
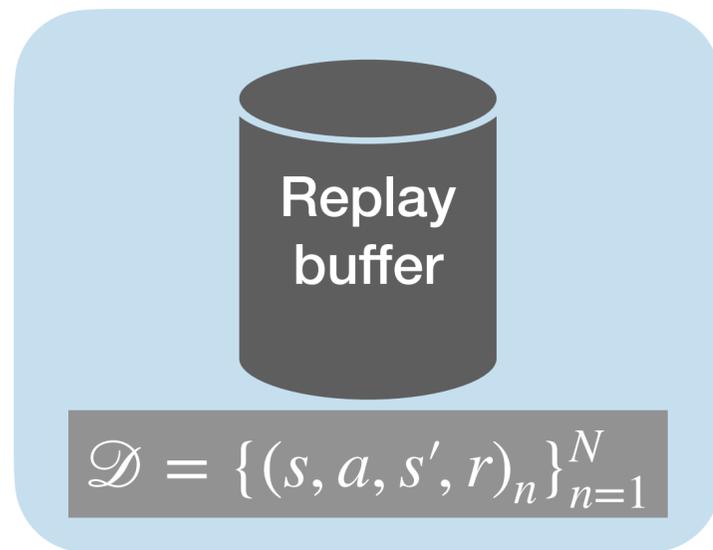
**Action-value function (aka Q-function):**

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi, P} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a, \pi \right]$$

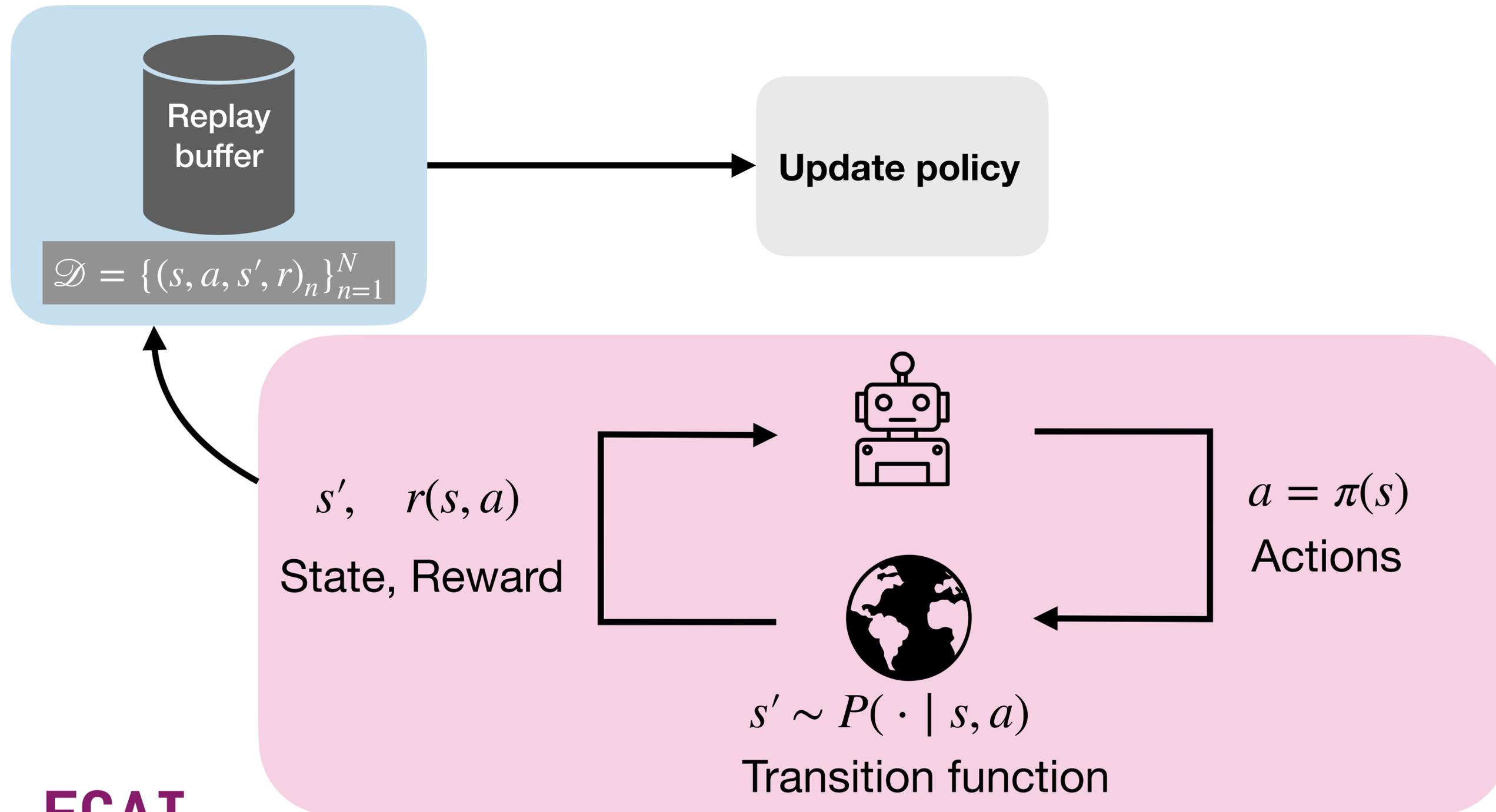
# Reinforcement Learning



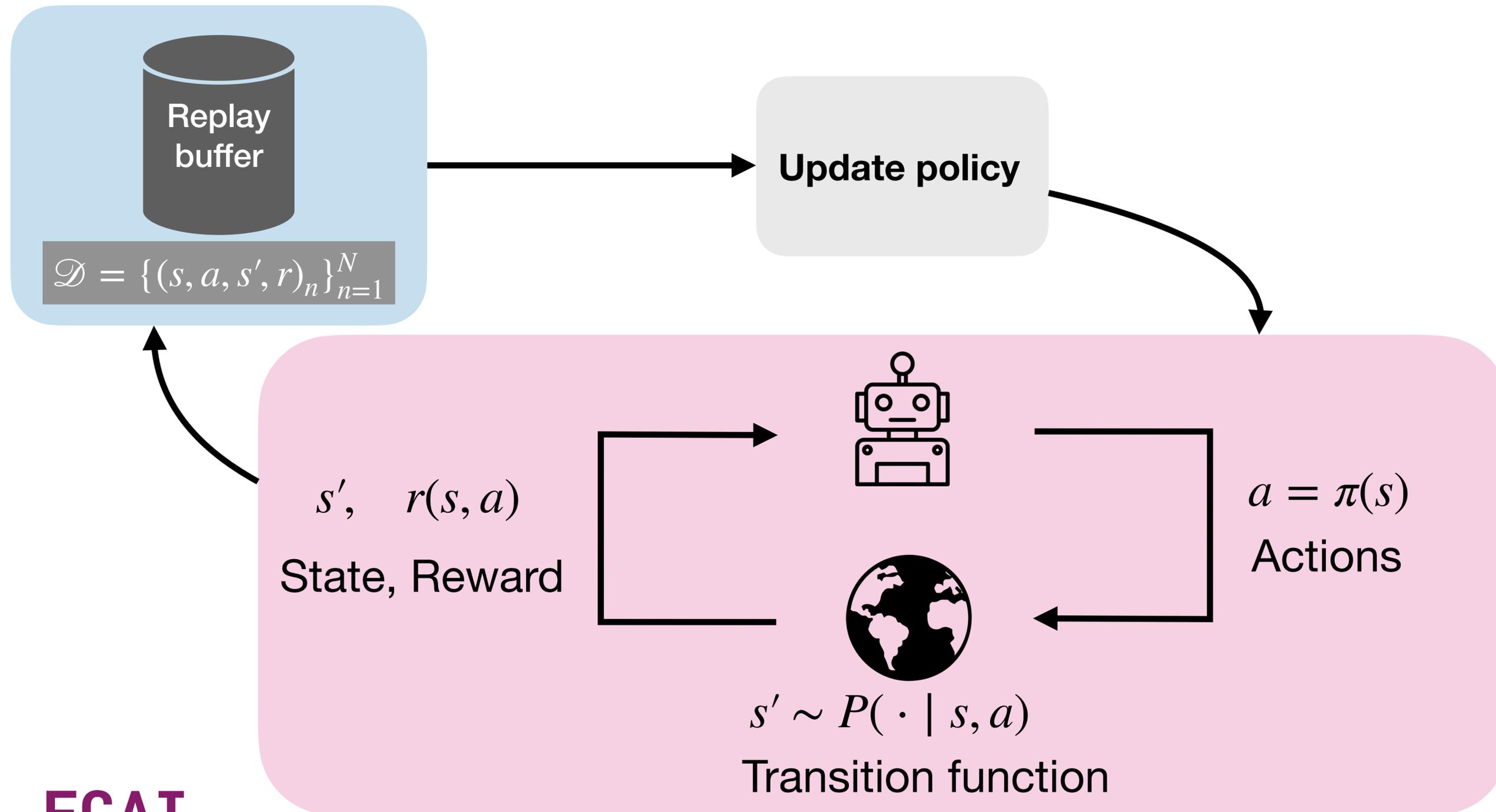
# Reinforcement Learning



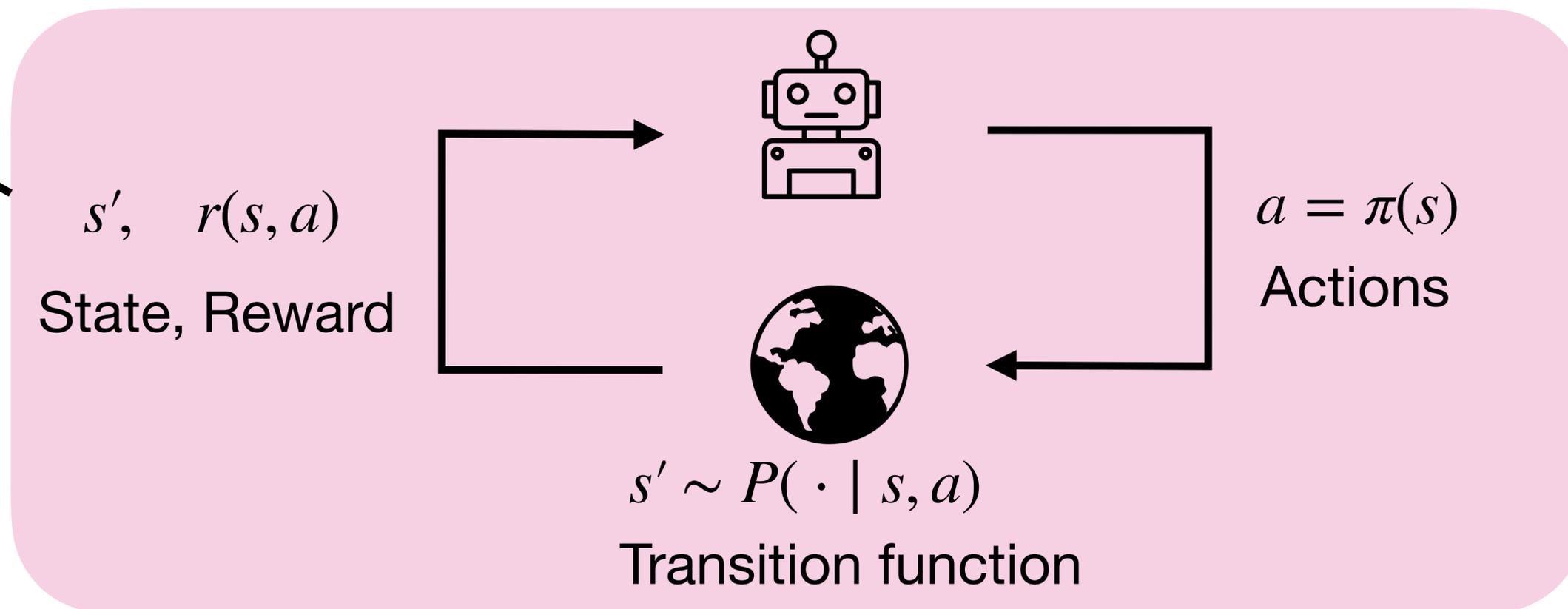
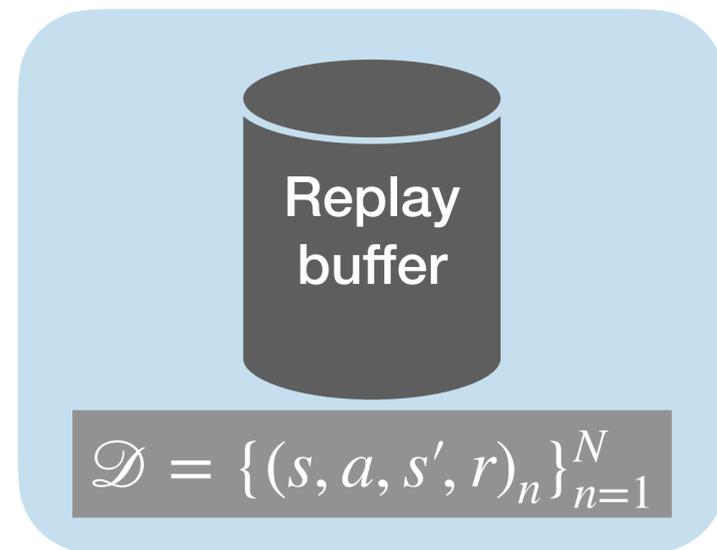
# Reinforcement Learning



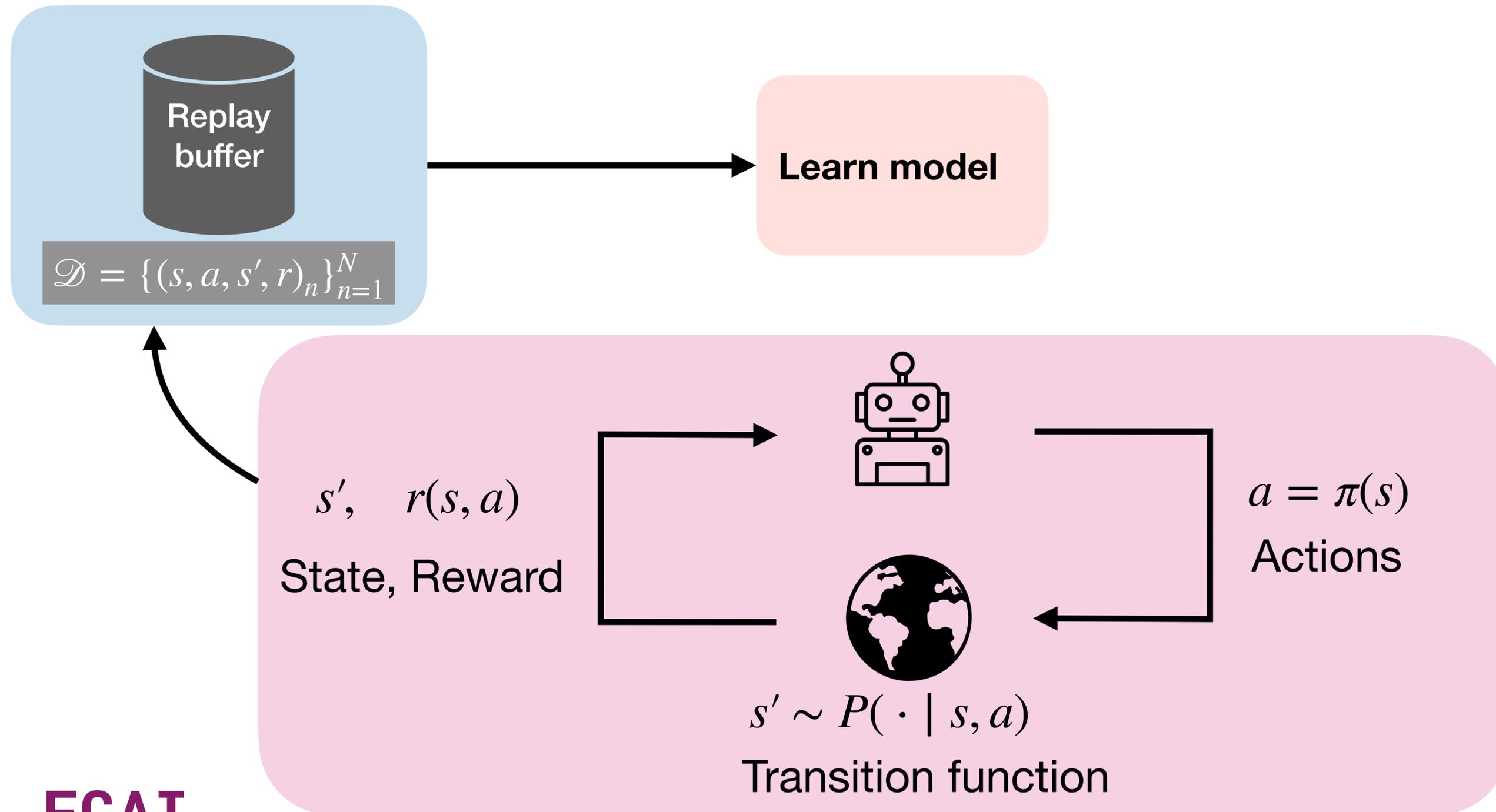
# Reinforcement Learning



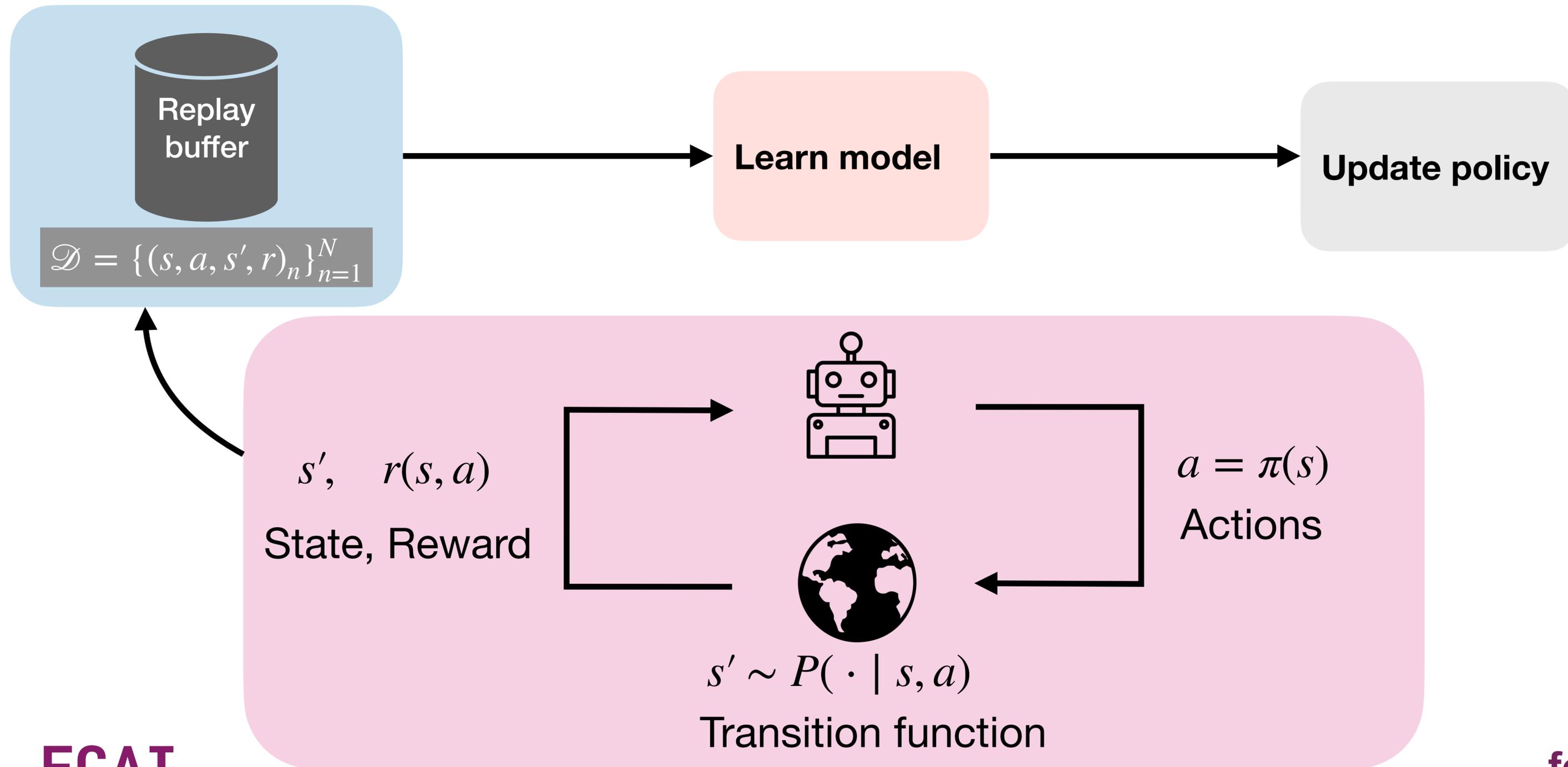
# Model-based Reinforcement Learning



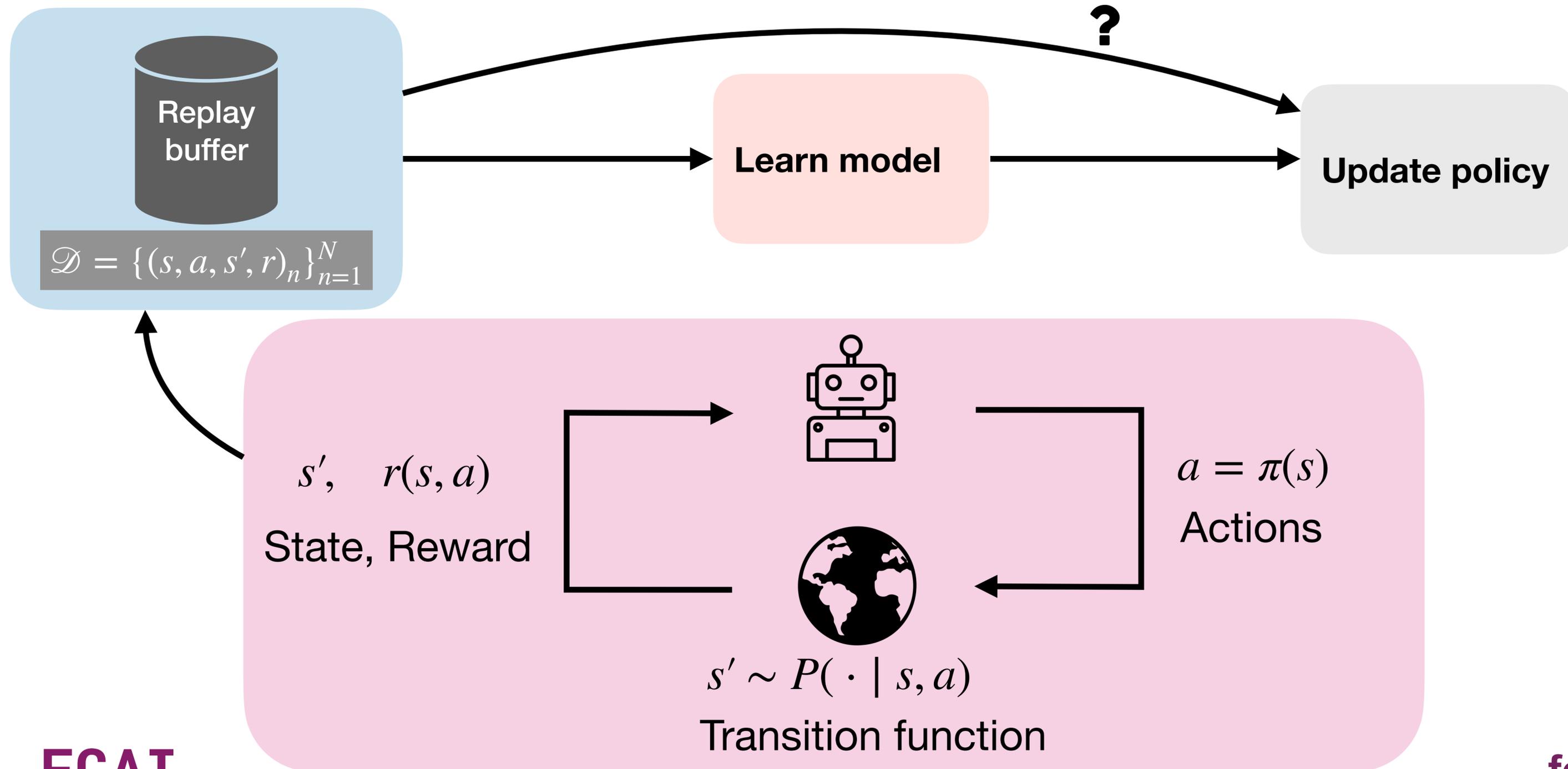
# Model-based Reinforcement Learning



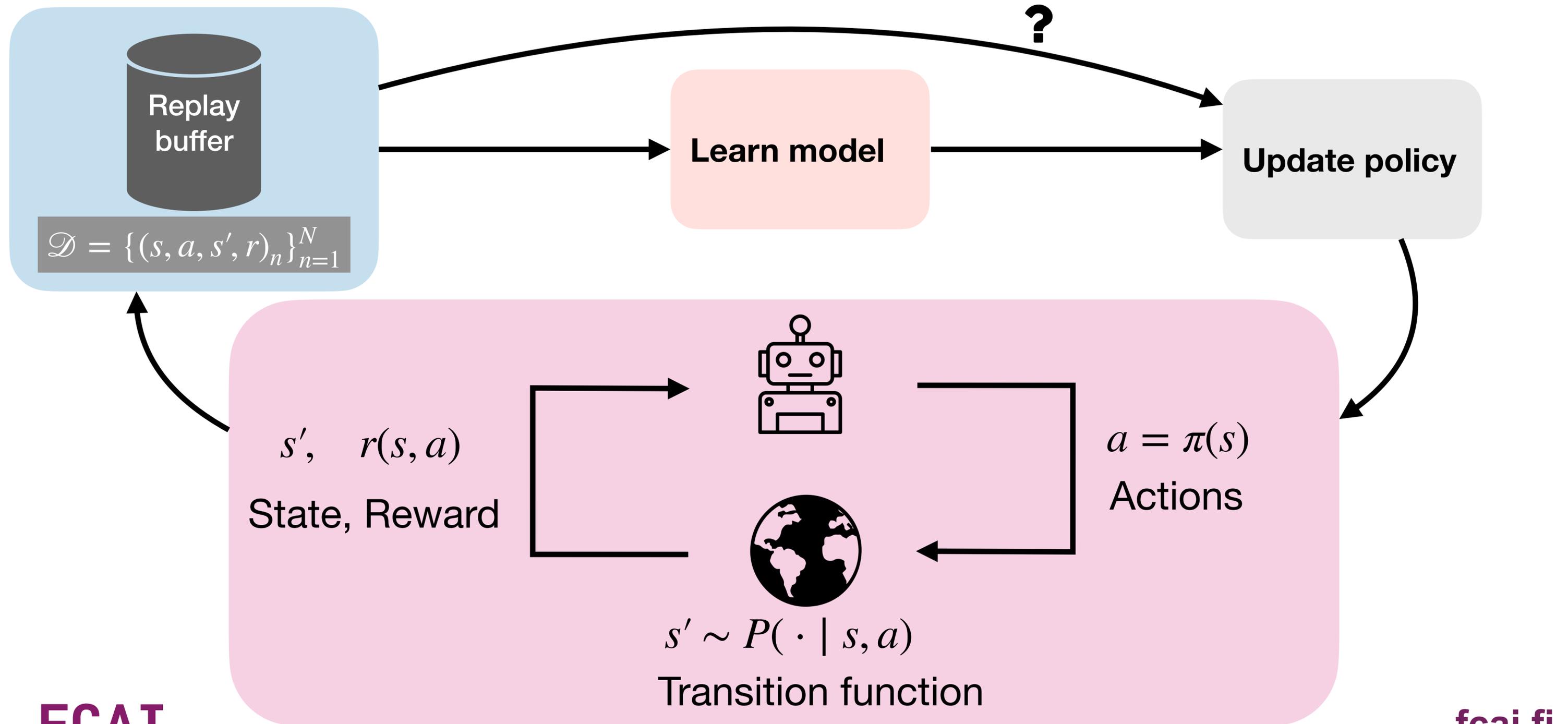
# Model-based Reinforcement Learning



# Model-based Reinforcement Learning



# Model-based Reinforcement Learning

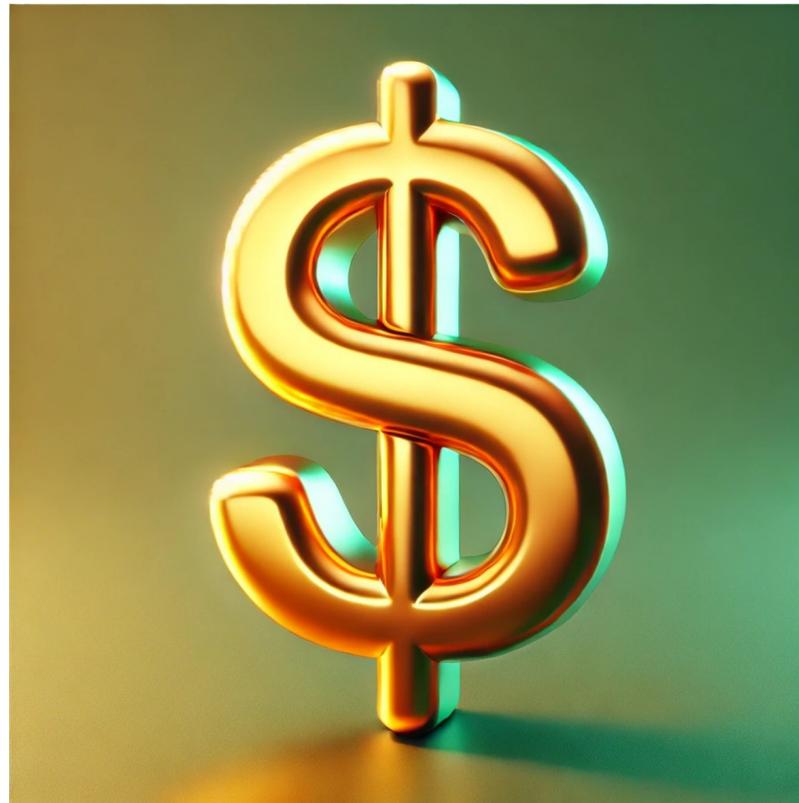


# Reinforcement Learning Has Its Drawbacks

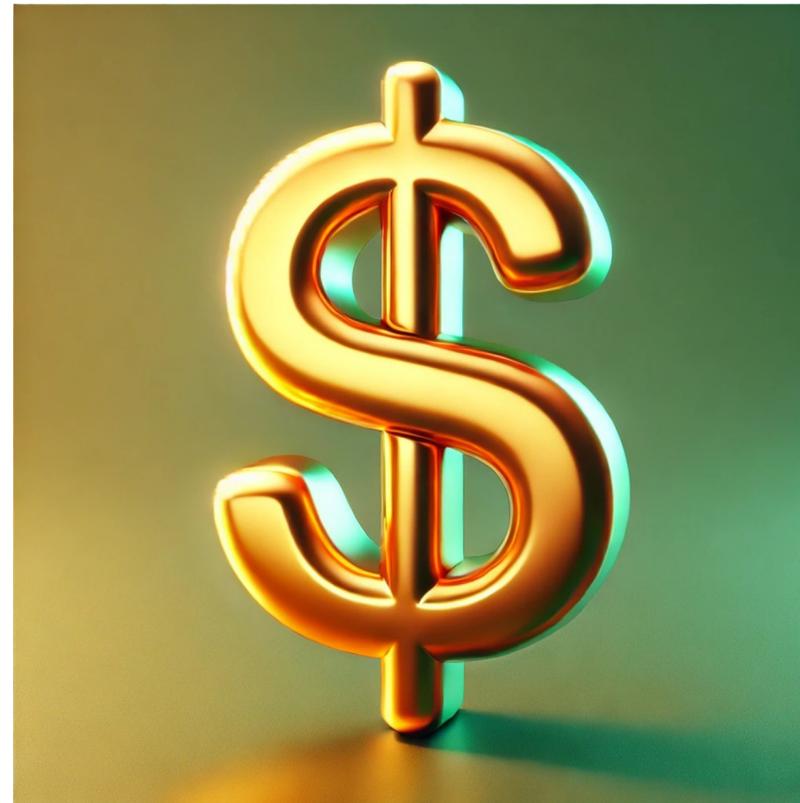
# Reinforcement Learning Has Its Drawbacks



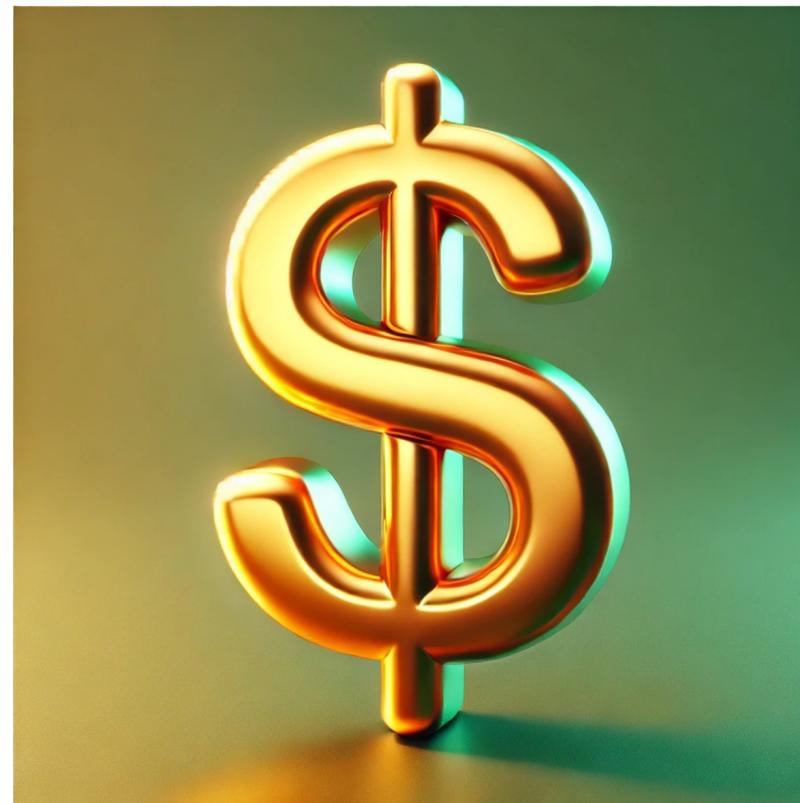
# Reinforcement Learning Has Its Drawbacks



# Reinforcement Learning Has Its Drawbacks



# Reinforcement Learning Has Its Drawbacks



**RL has a sample efficiency problem!**

# Model-free vs Model-based RL

# Model-free vs Model-based RL

	<b>Model-free</b>	<b>Model-based</b>
--	-------------------	--------------------

# Model-free vs Model-based RL

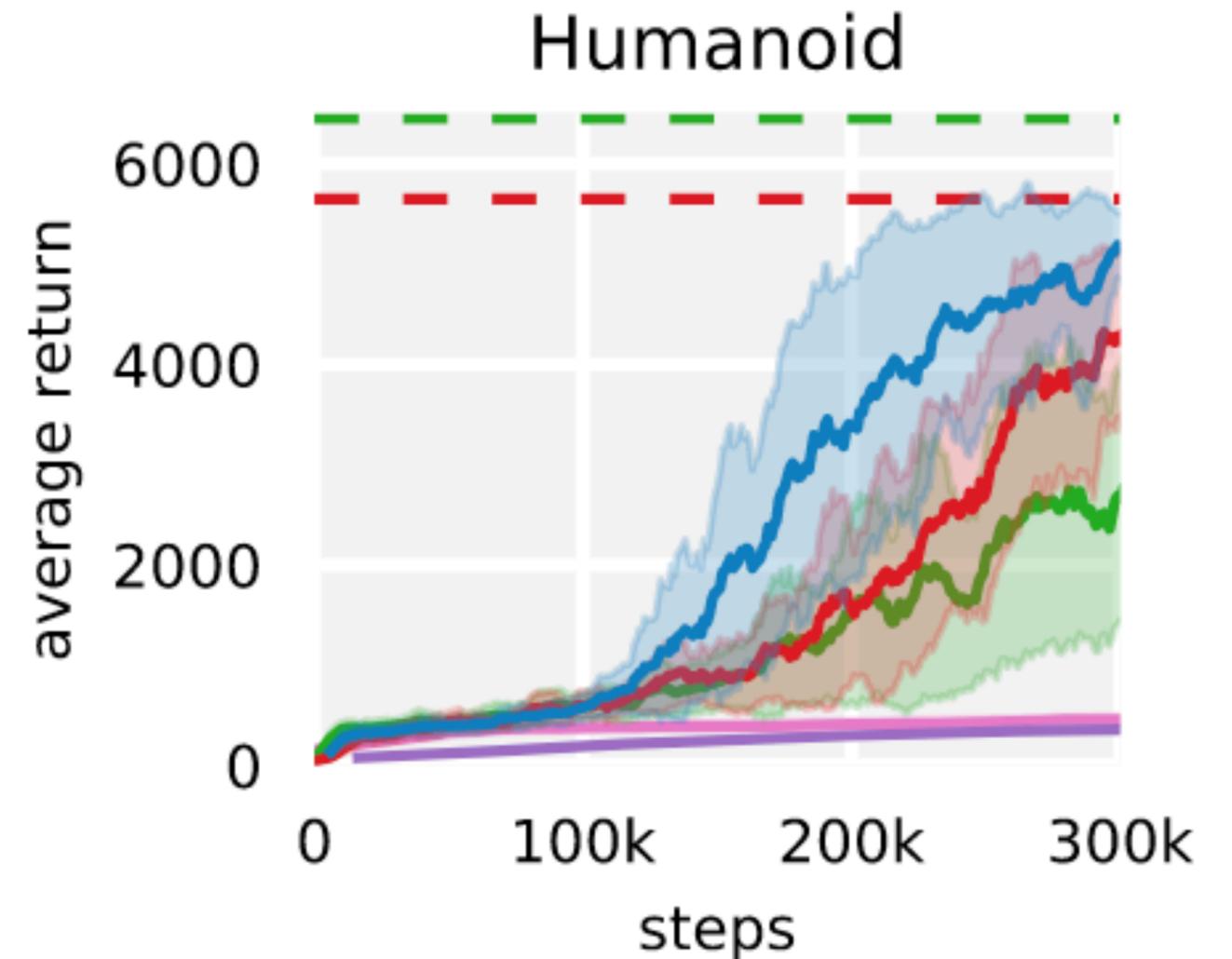
	Model-free	Model-based
Asymptotic performance	✓	Depends

# Model-free vs Model-based RL

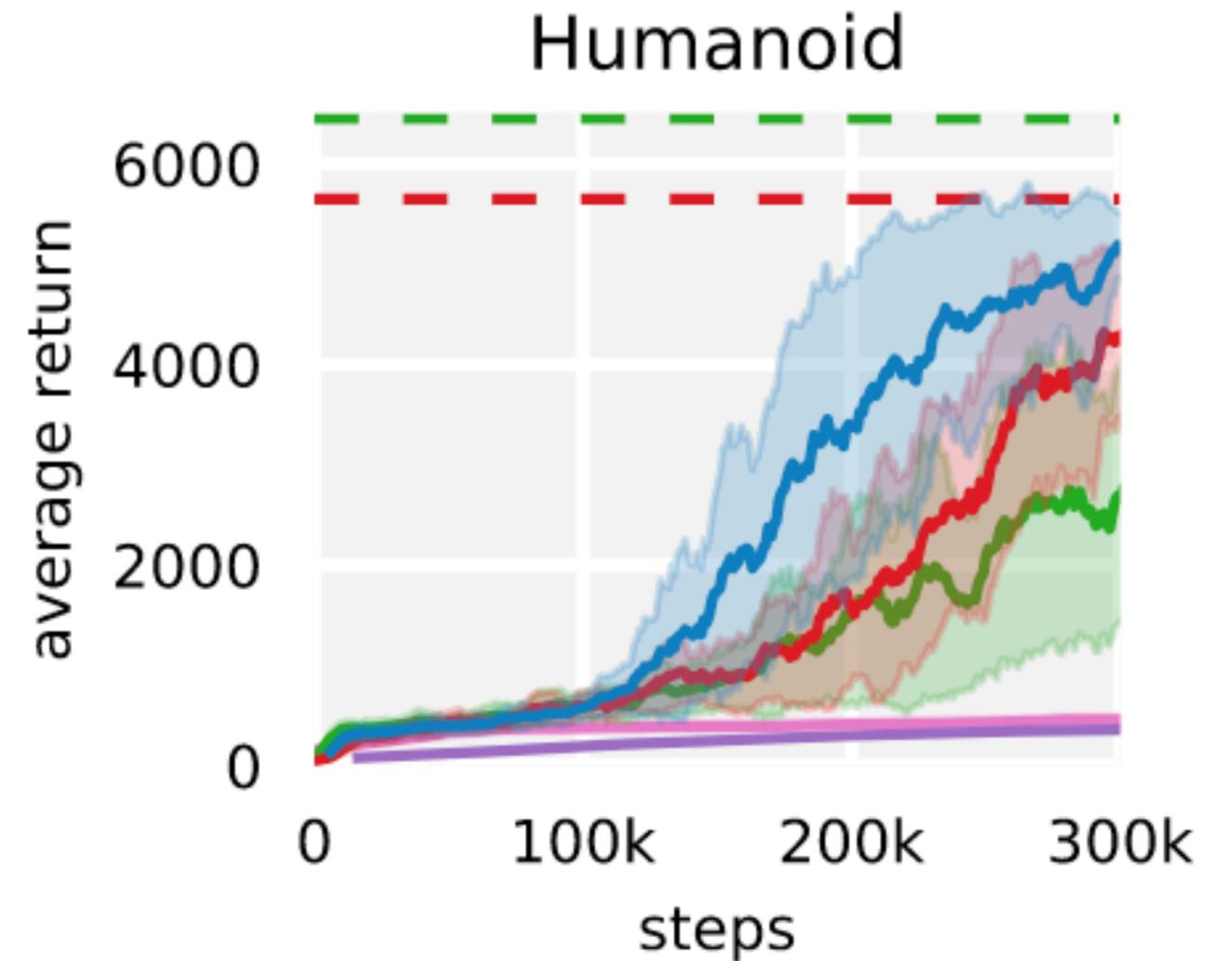
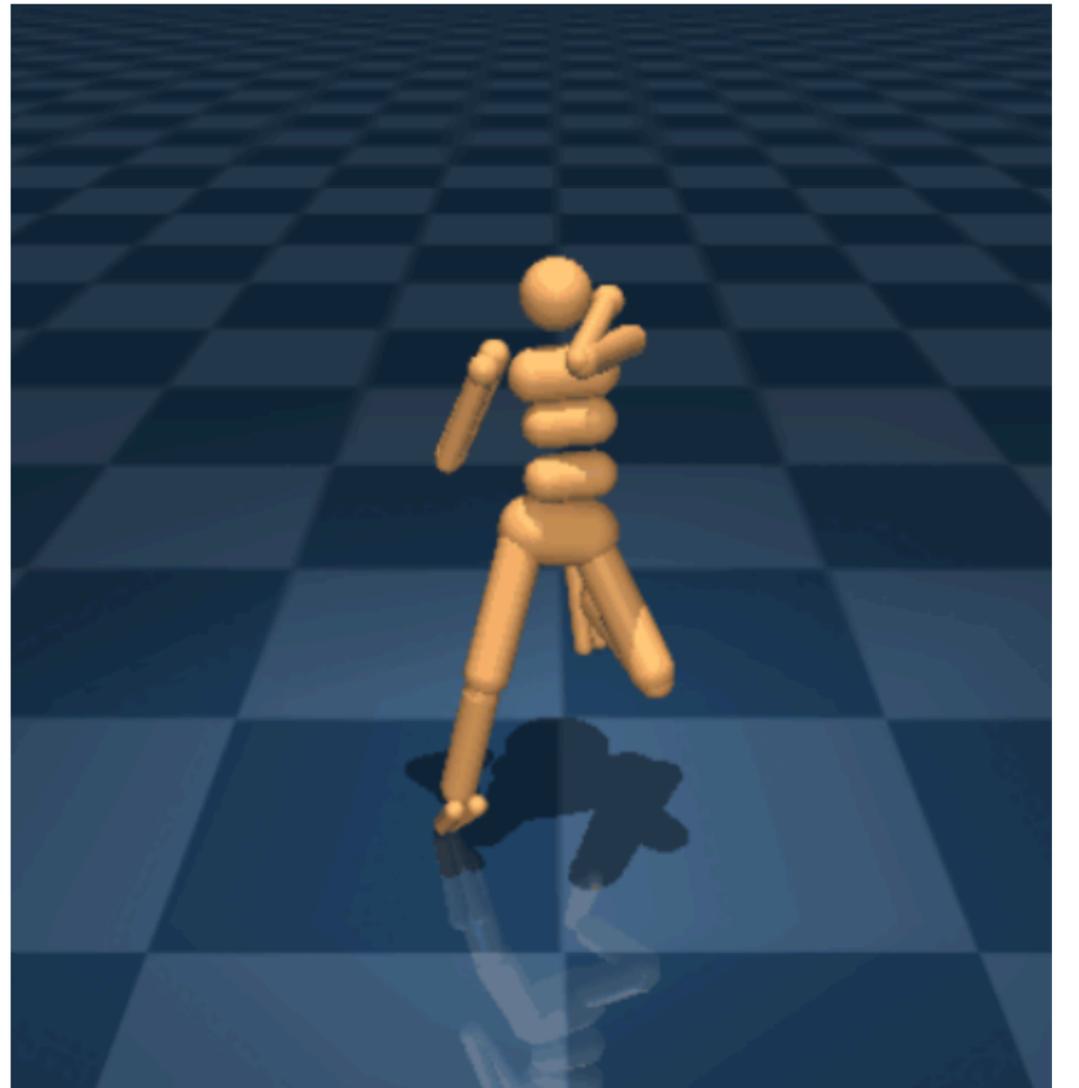
	Model-free	Model-based
Asymptotic performance	✓	Depends
Sample efficiency	✗	✓

# Model-free vs Model-based RL

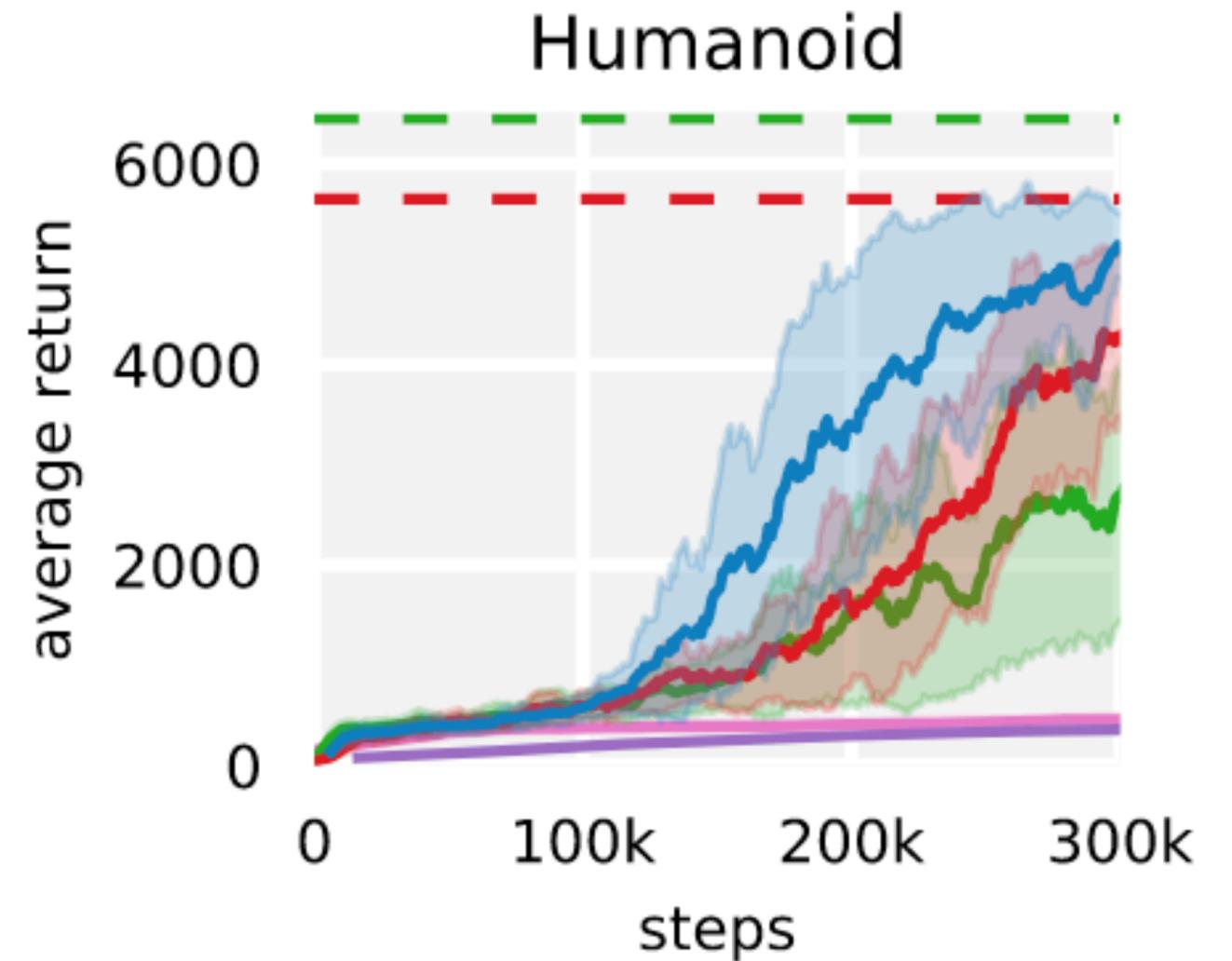
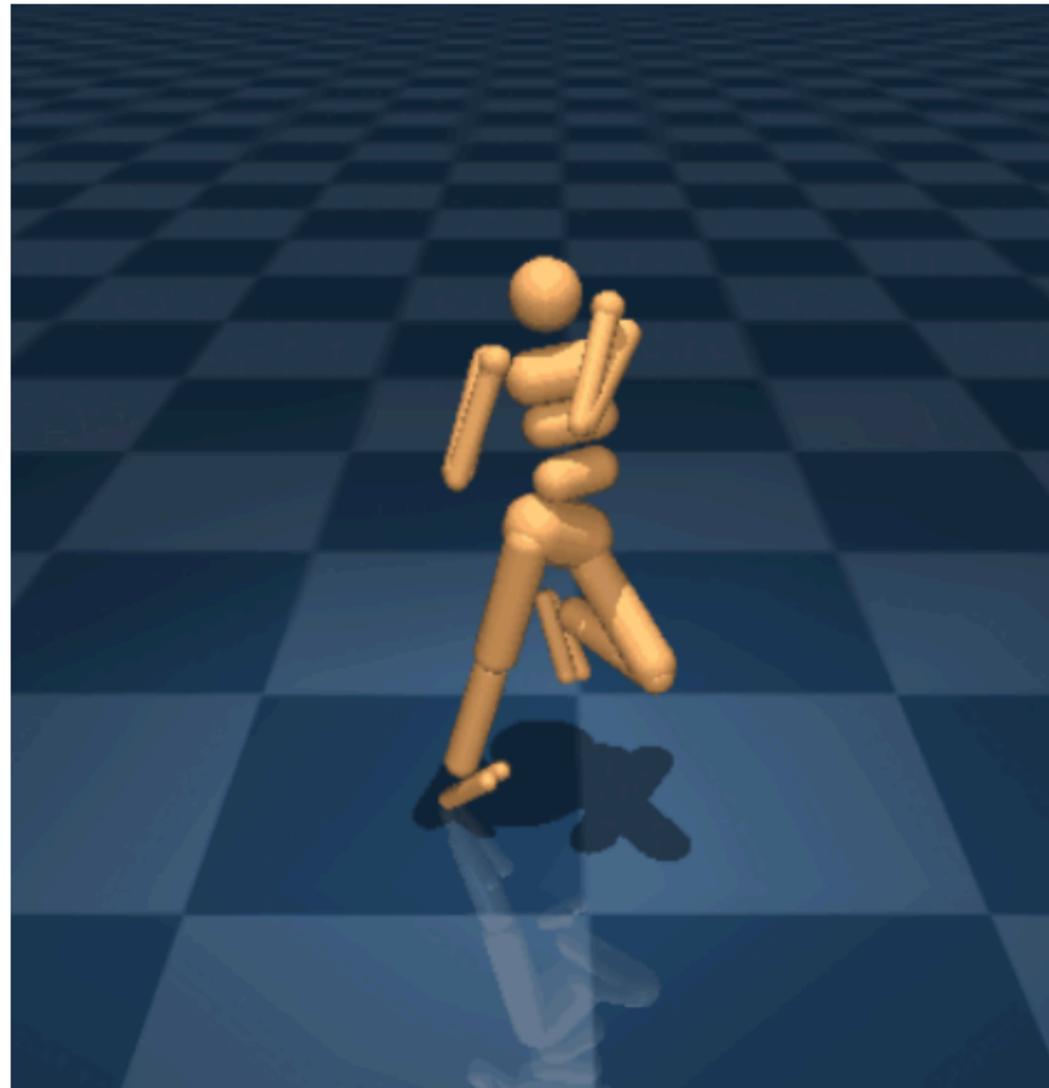
	Model-
Asymptotic performance	✓
Sample efficiency	✗



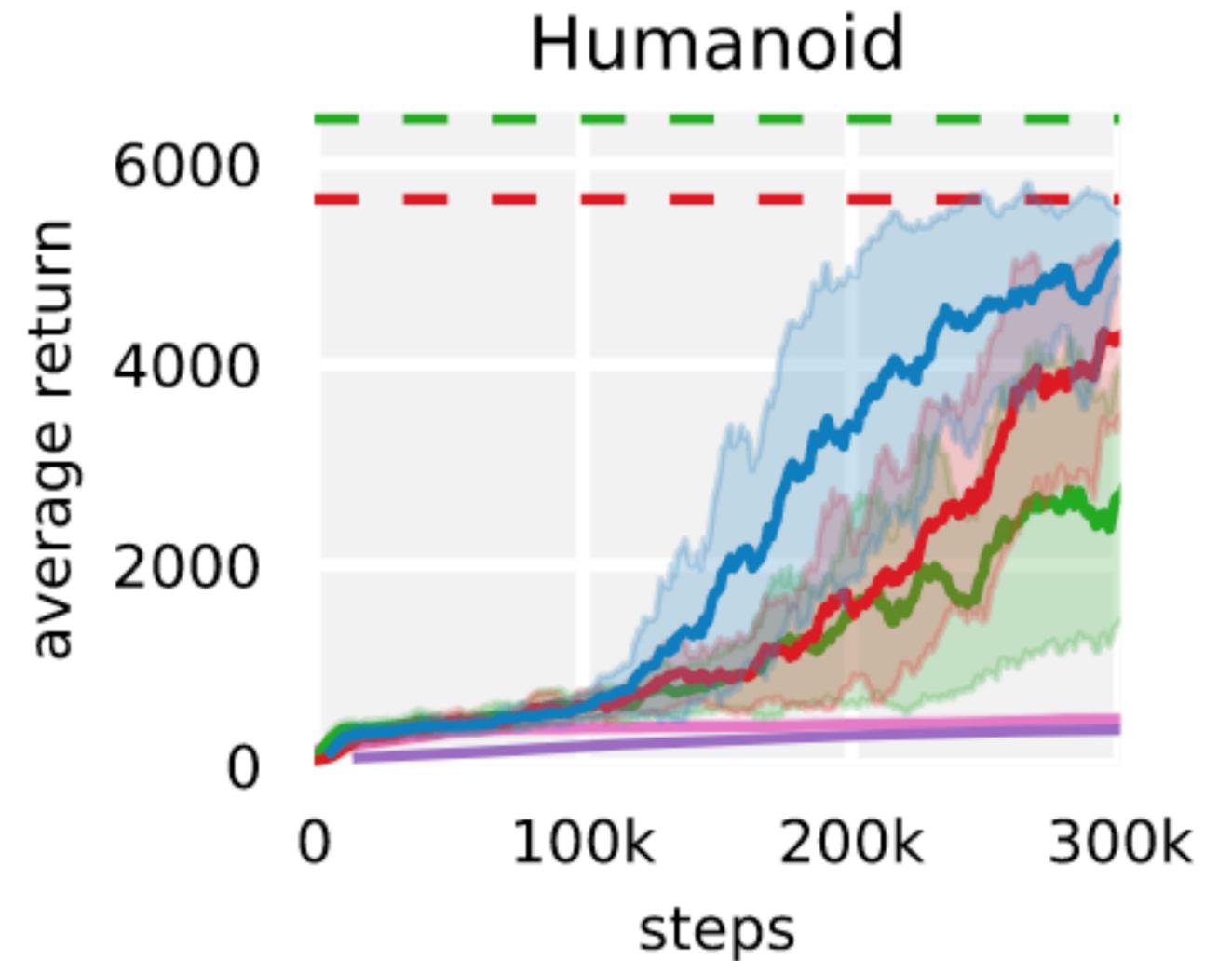
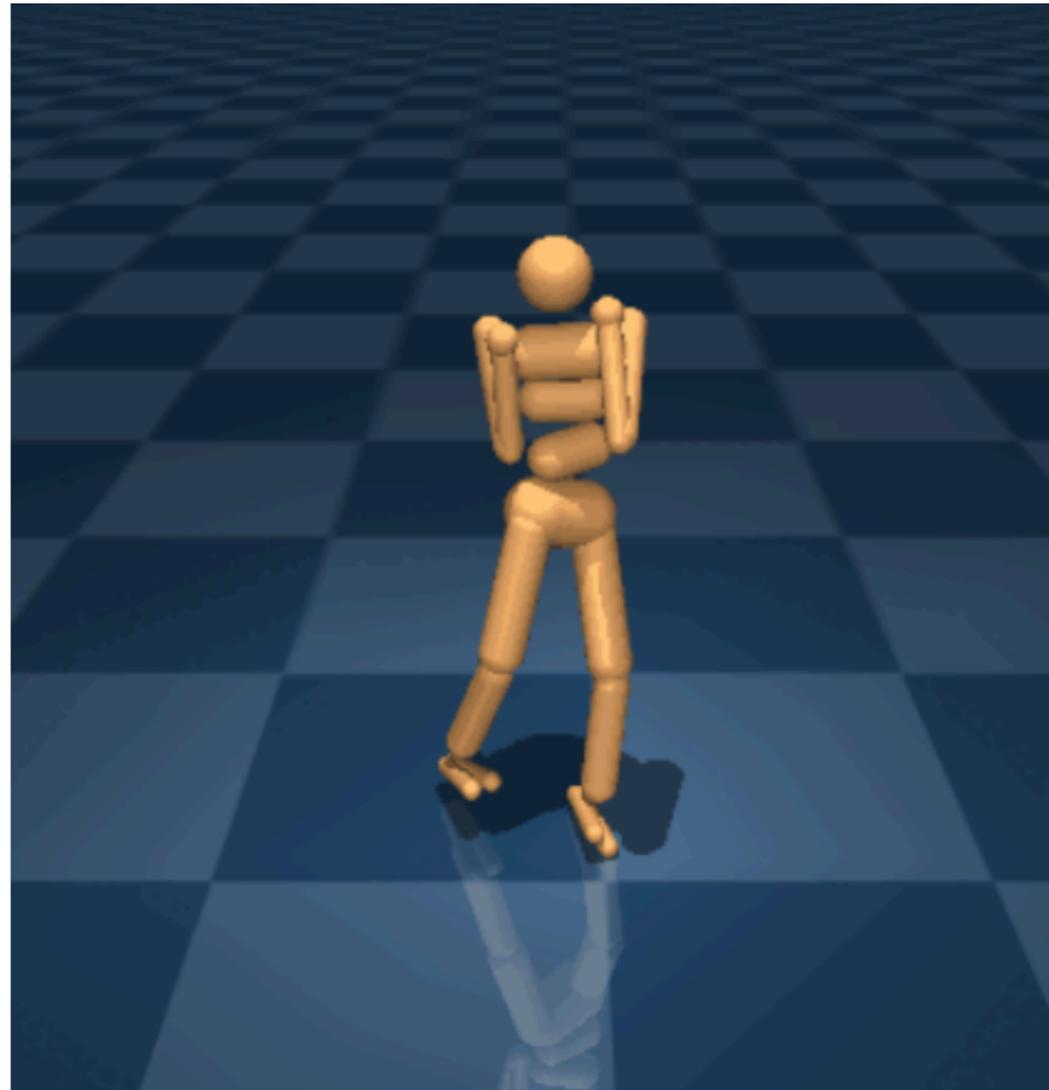
# Model-free vs Model-based RL



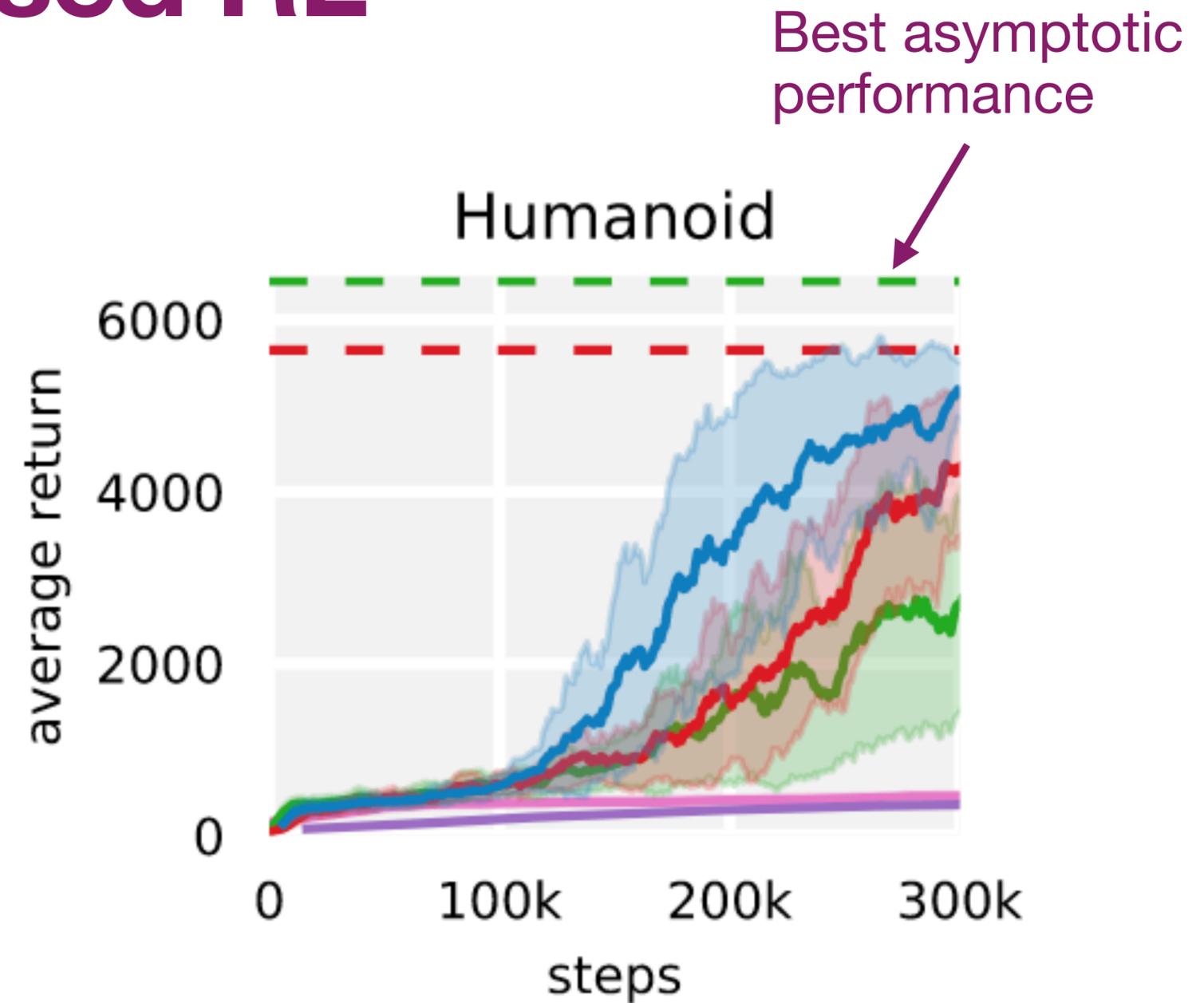
# Model-free vs Model-based RL



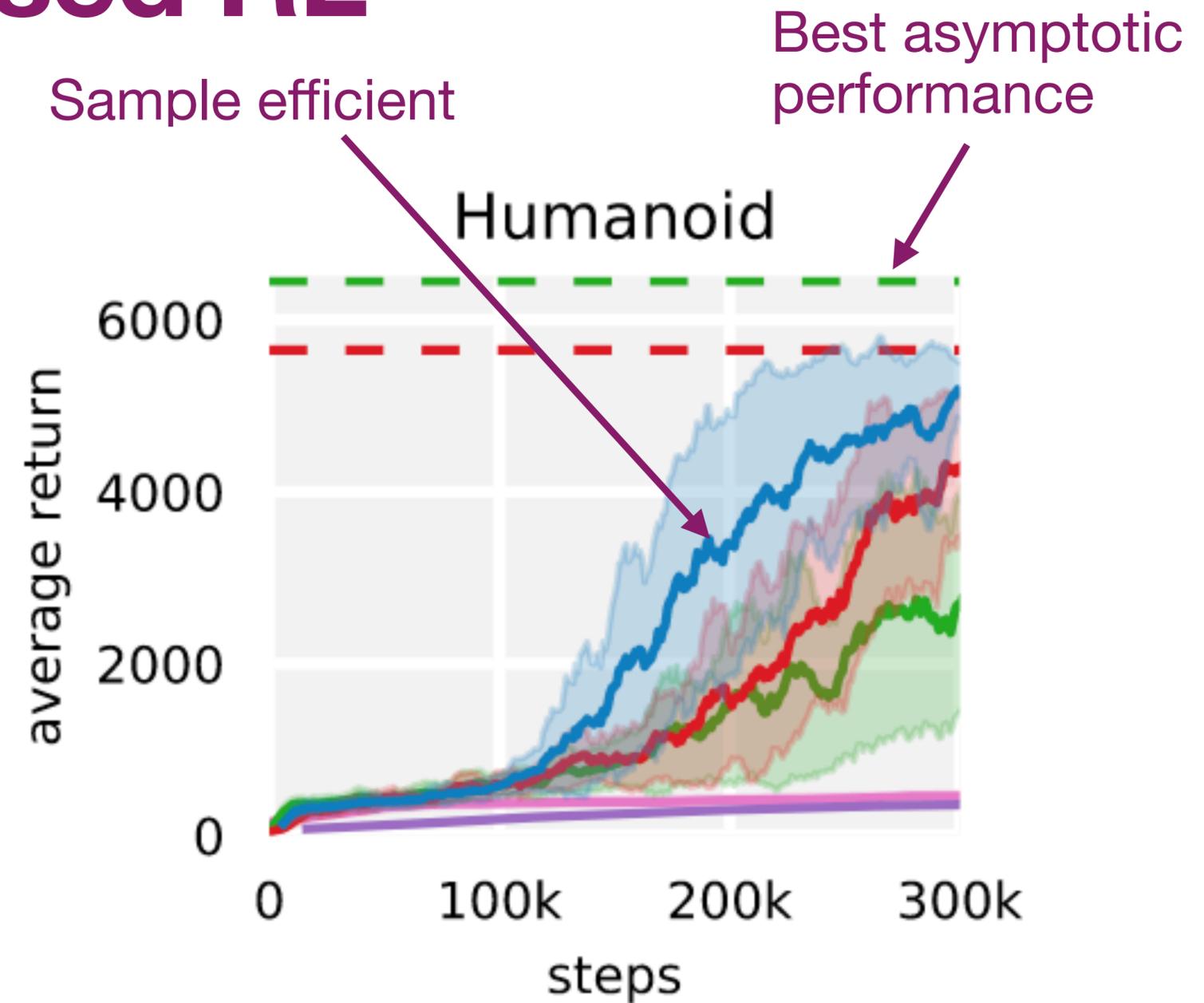
# Model-free vs Model-based RL



# Model-free vs Model-based RL



# Model-free vs Model-based RL



# Model-free vs Model-based RL

	Model-free	Model-based
Asymptotic performance	✓	Depends
Sample efficiency	✗	✓

# Model-free vs Model-based RL

	Model-free	Model-based
Asymptotic performance	✓	Depends
Sample efficiency	✗	✓
Computation at deployment	✓	✗/✓

# Model-free vs Model-based RL

	Model-free	Model-based
Asymptotic performance	✓	Depends
Sample efficiency	✗	✓
Computation at deployment	✓	✗/✓
Adapting to new tasks	✗	✓

# Model-free vs Model-based RL

	Model-free	Model-based
Asymptotic performance	✓	Depends
Sample efficiency	✗	✓
Computation at deployment	✓	✗/✓
Adapting to new tasks	✗	✓
Exploration	✗	✓

# What is a “Model”?

*Definition: a model is a representation that **explicitly** encodes knowledge about the structure of the environment and task.*

# What is a “Model”?

*Definition: a model is a representation that **explicitly** encodes knowledge about the structure of the environment and task.*

Dynamics (transition) model

$$s_{t+1} = f(s_t, a_t)$$

# What is a “Model”?

*Definition: a model is a representation that **explicitly** encodes knowledge about the structure of the environment and task.*

Dynamics (transition) model

$$s_{t+1} = f(s_t, a_t)$$

Reward model

$$r_{t+1} = f(s_t, a_t)$$

# What is a “Model”?

*Definition: a model is a representation that **explicitly** encodes knowledge about the structure of the environment and task.*

Dynamics (transition) model

$$s_{t+1} = f(s_t, a_t)$$

Reward model

$$r_{t+1} = f(s_t, a_t)$$

Inverse dynamics model

$$a_t = f^{-1}(s_t, s_{t+1})$$

# What is a “Model”?

*Definition: a model is a representation that **explicitly** encodes knowledge about the structure of the environment and task.*

Dynamics (transition) model

$$s_{t+1} = f(s_t, a_t)$$

Reward model

$$r_{t+1} = f(s_t, a_t)$$

Inverse dynamics model

$$a_t = f^{-1}(s_t, s_{t+1})$$

Model of distance

$$d_{ij} = f_d(s_i, s_j)$$

# What is a “Model”?

*Definition: a model is a representation that **explicitly** encodes knowledge about the structure of the environment and task.*

Dynamics (transition) model

$$s_{t+1} = f(s_t, a_t)$$

Reward model

$$r_{t+1} = f(s_t, a_t)$$

Inverse dynamics model

$$a_t = f^{-1}(s_t, s_{t+1})$$

Model of distance

$$d_{ij} = f_d(s_i, s_j)$$

Model of future returns

$$G_t = Q(s_t, a_t)$$

# What is a “Model”?

*Definition: a model is a representation that **explicitly** encodes knowledge about the structure of the environment and task.*

Dynamics (transition) model

$$s_{t+1} = f(s_t, a_t)$$

Reward model

$$r_{t+1} = f(s_t, a_t)$$

Typically this is what's meant in model-based RL

Inverse dynamics model

$$a_t = f^{-1}(s_t, s_{t+1})$$

Model of distance

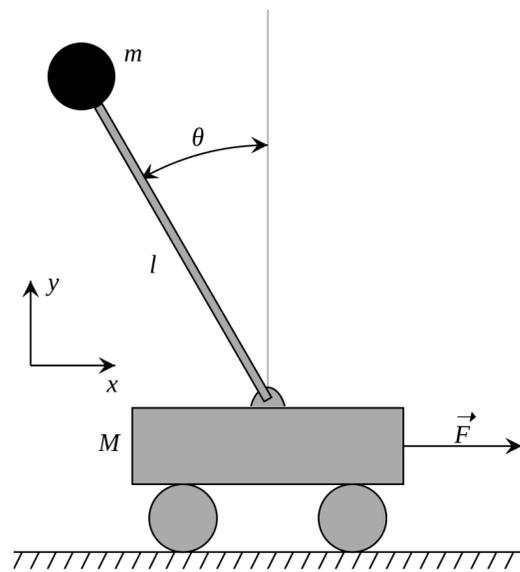
$$d_{ij} = f_d(s_i, s_j)$$

Model of future returns

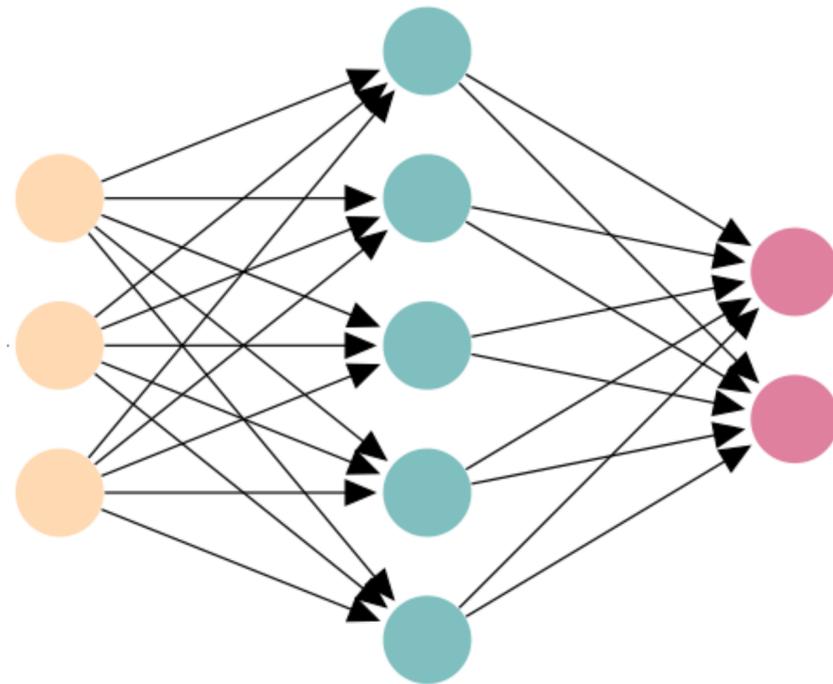
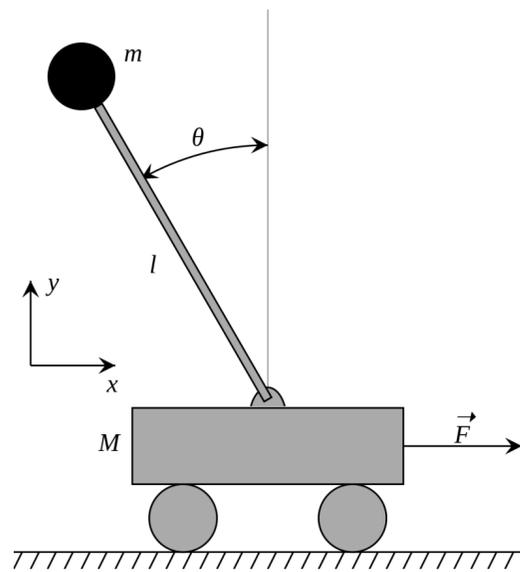
$$G_t = Q(s_t, a_t)$$

# What is a “Model”?

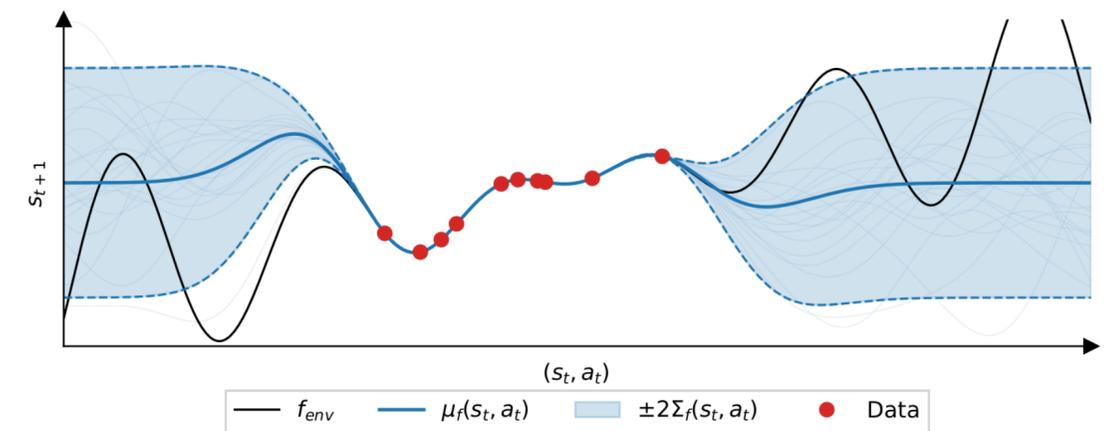
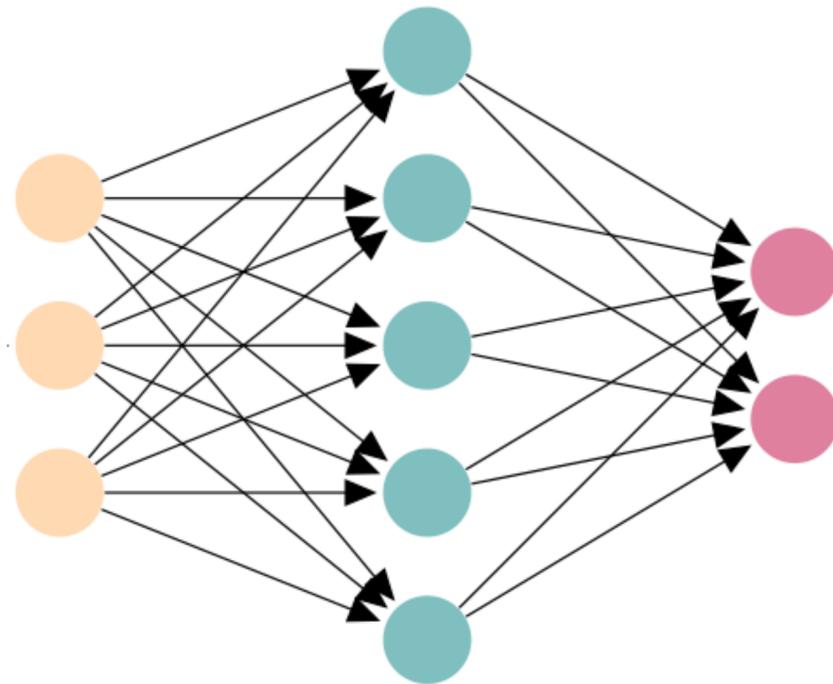
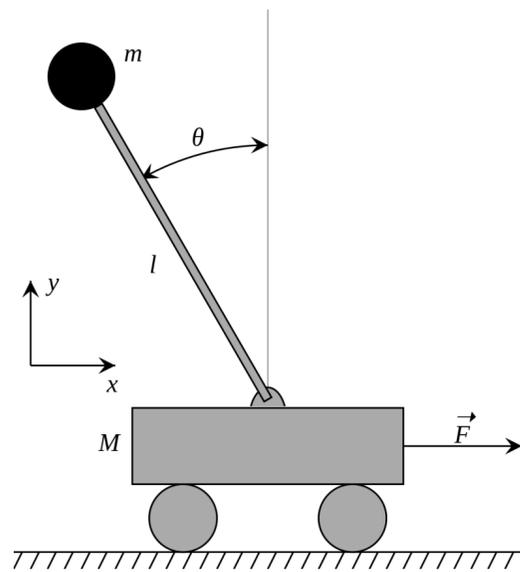
# What is a “Model”?



# What is a “Model”?



# What is a “Model”?



# Learning Objectives

Understand

1. ~~What a “model” is in model-based RL~~
2. How a “model” can aid decision making
3. Differences between background and decision-time planning
4. Decision-time planning strategies for continuous actions
5. Sources of uncertainty in model-based RL
6. Rationale and insights for decision-making under uncertainty

# Planning

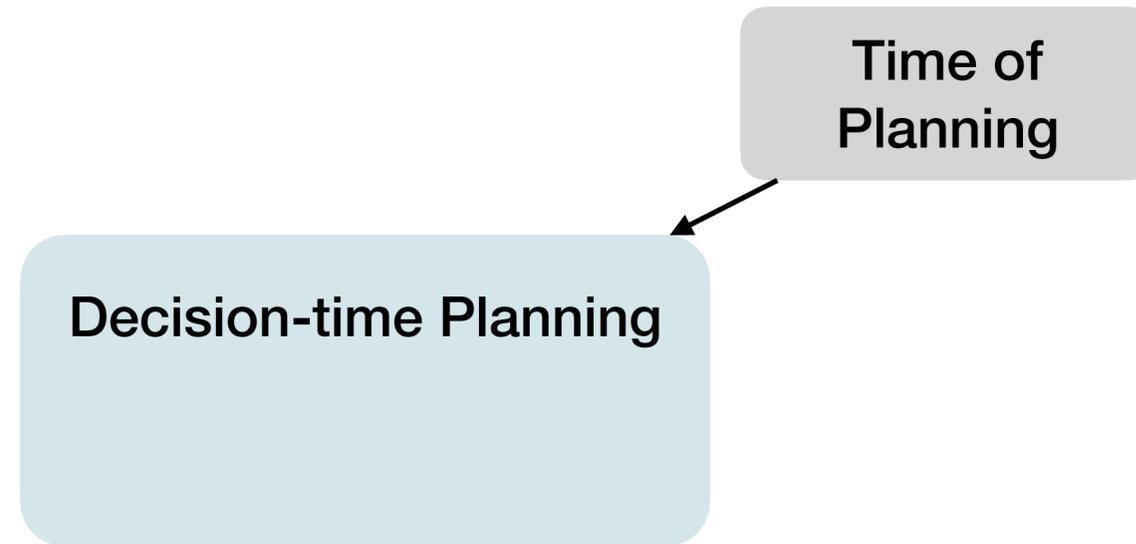
# How Do We Use The "Model"?

## Background vs Decision-time Planning

Time of  
Planning

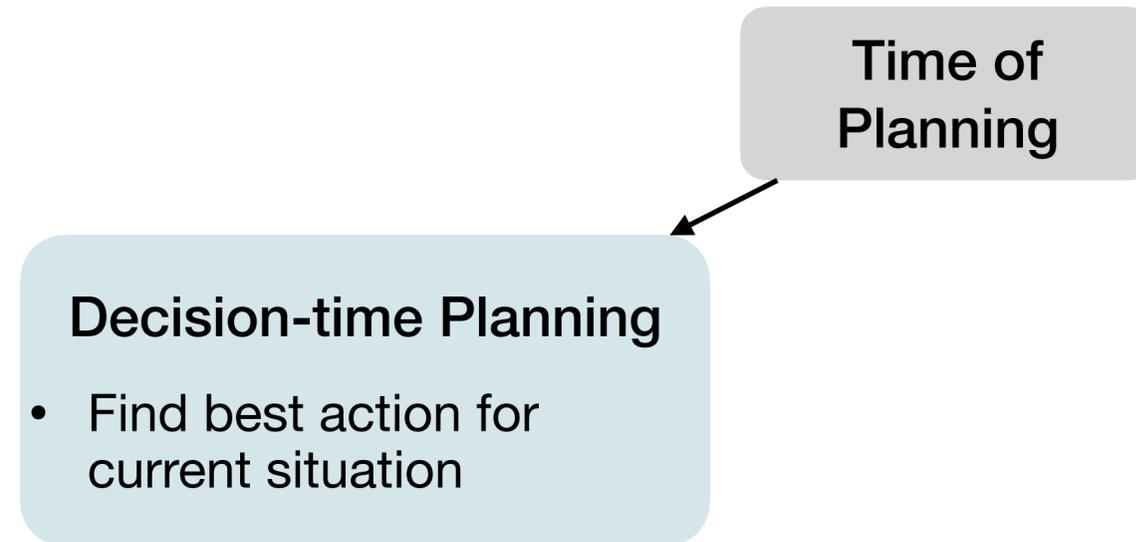
# How Do We Use The "Model"?

## Background vs Decision-time Planning



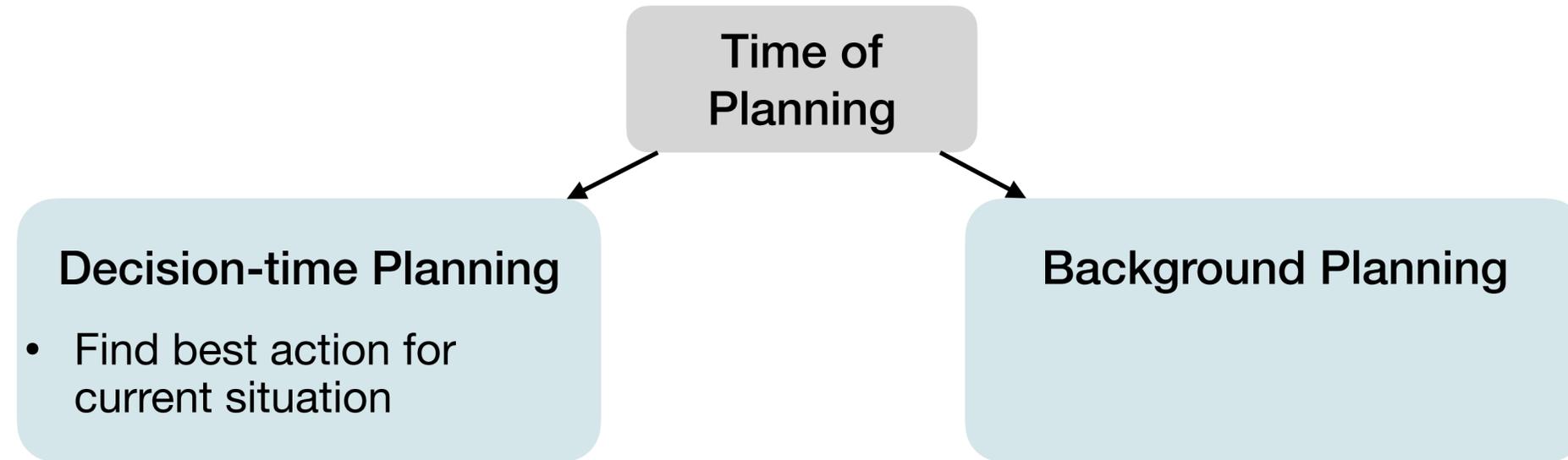
# How Do We Use The "Model"?

## Background vs Decision-time Planning



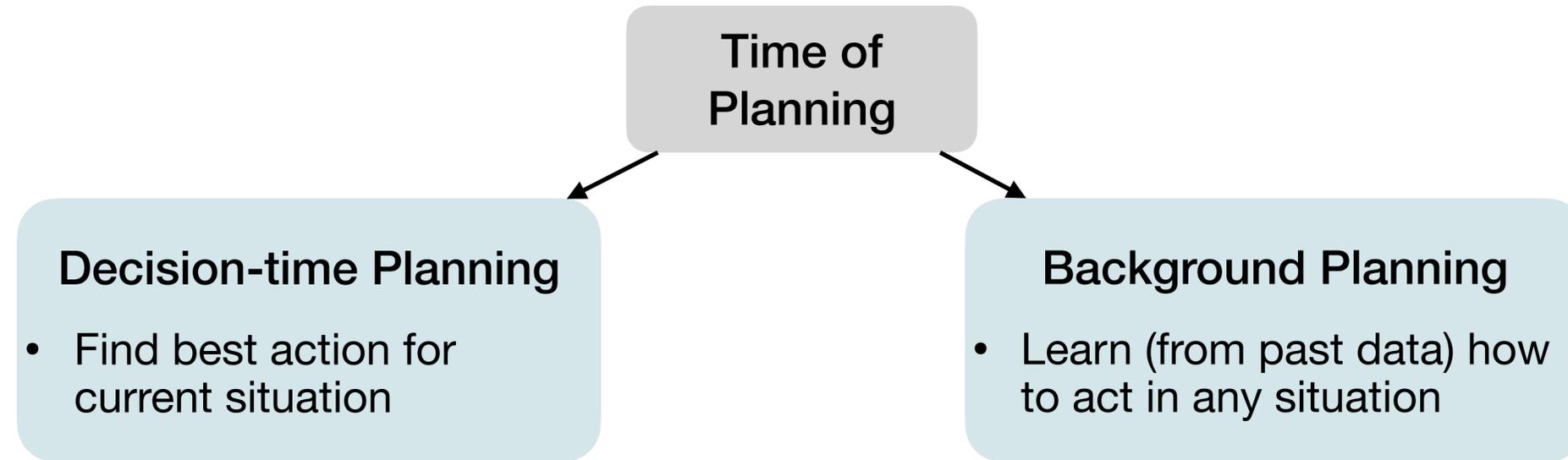
# How Do We Use The "Model"?

## Background vs Decision-time Planning



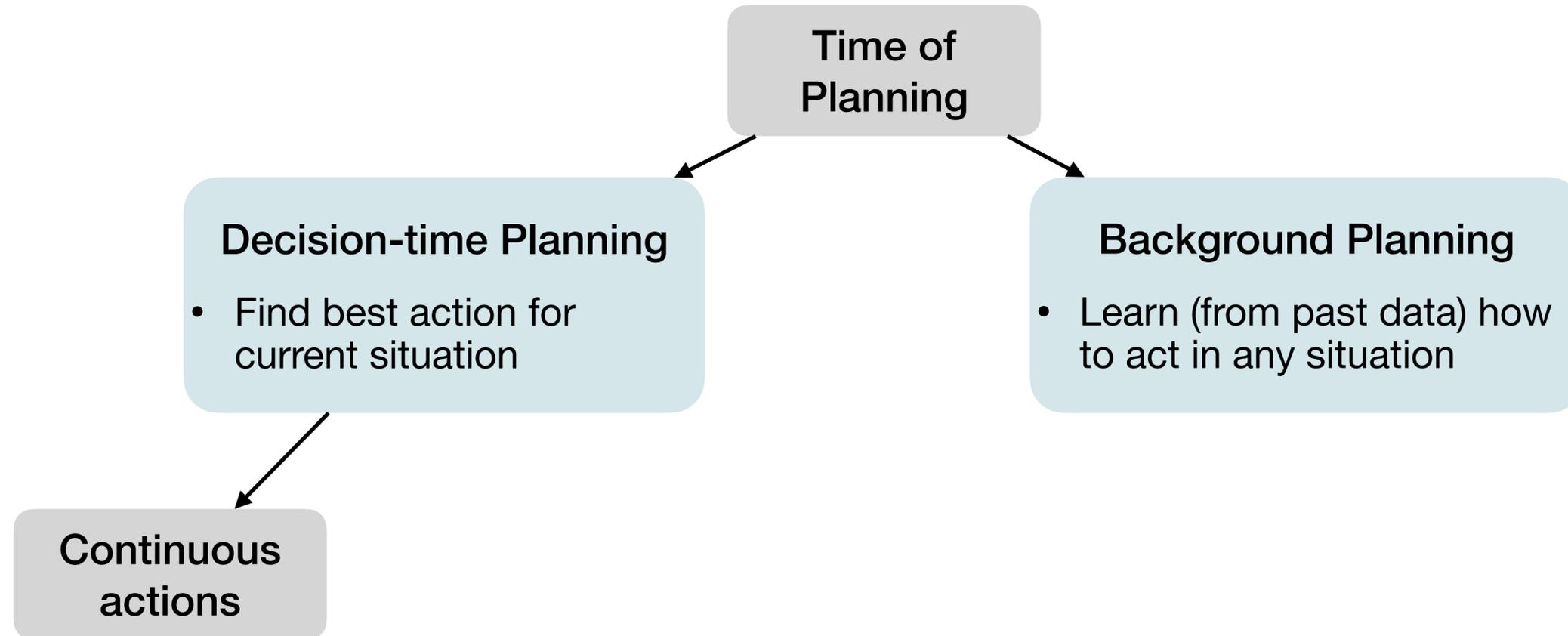
# How Do We Use The "Model"?

## Background vs Decision-time Planning



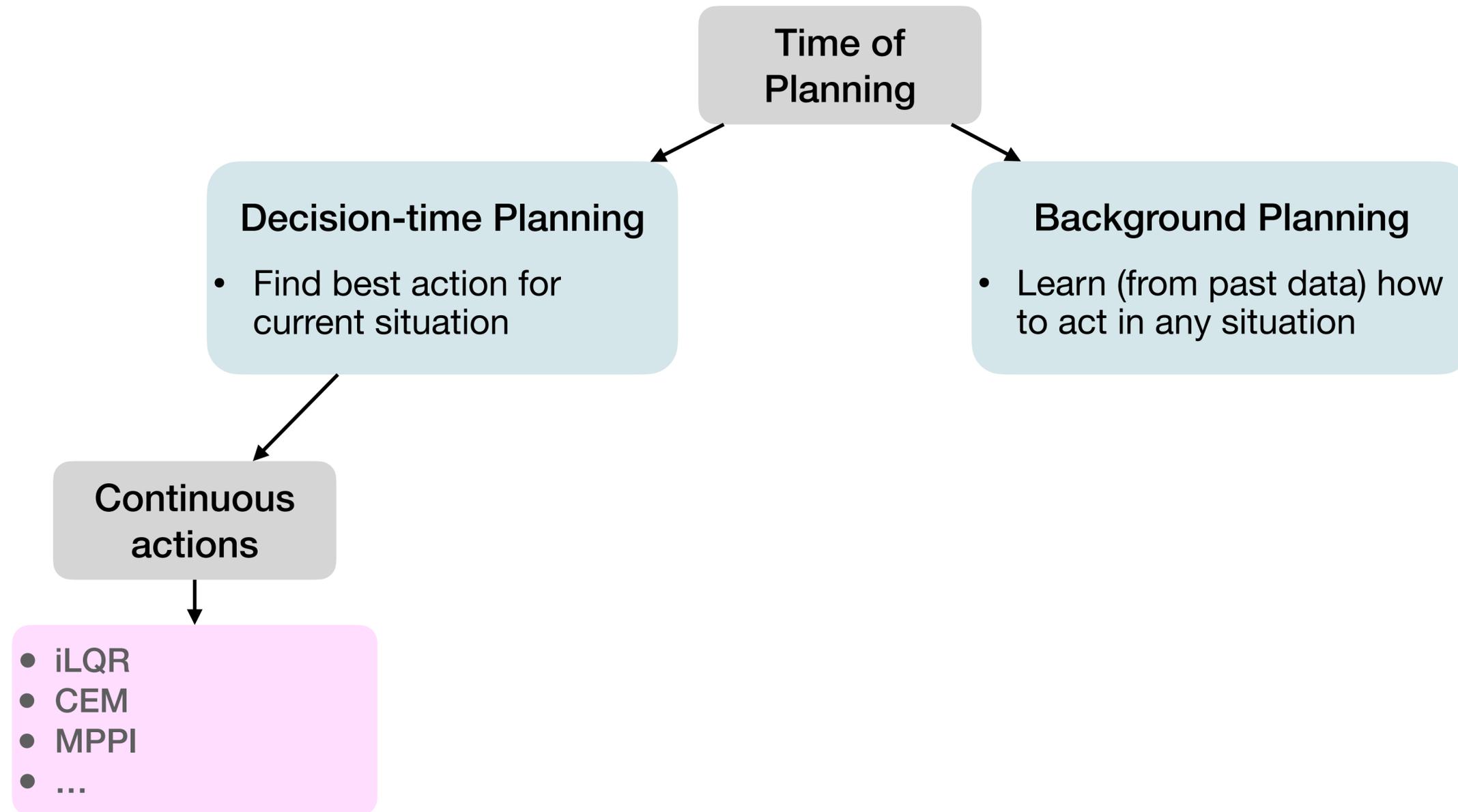
# How Do We Use The "Model"?

## Background vs Decision-time Planning



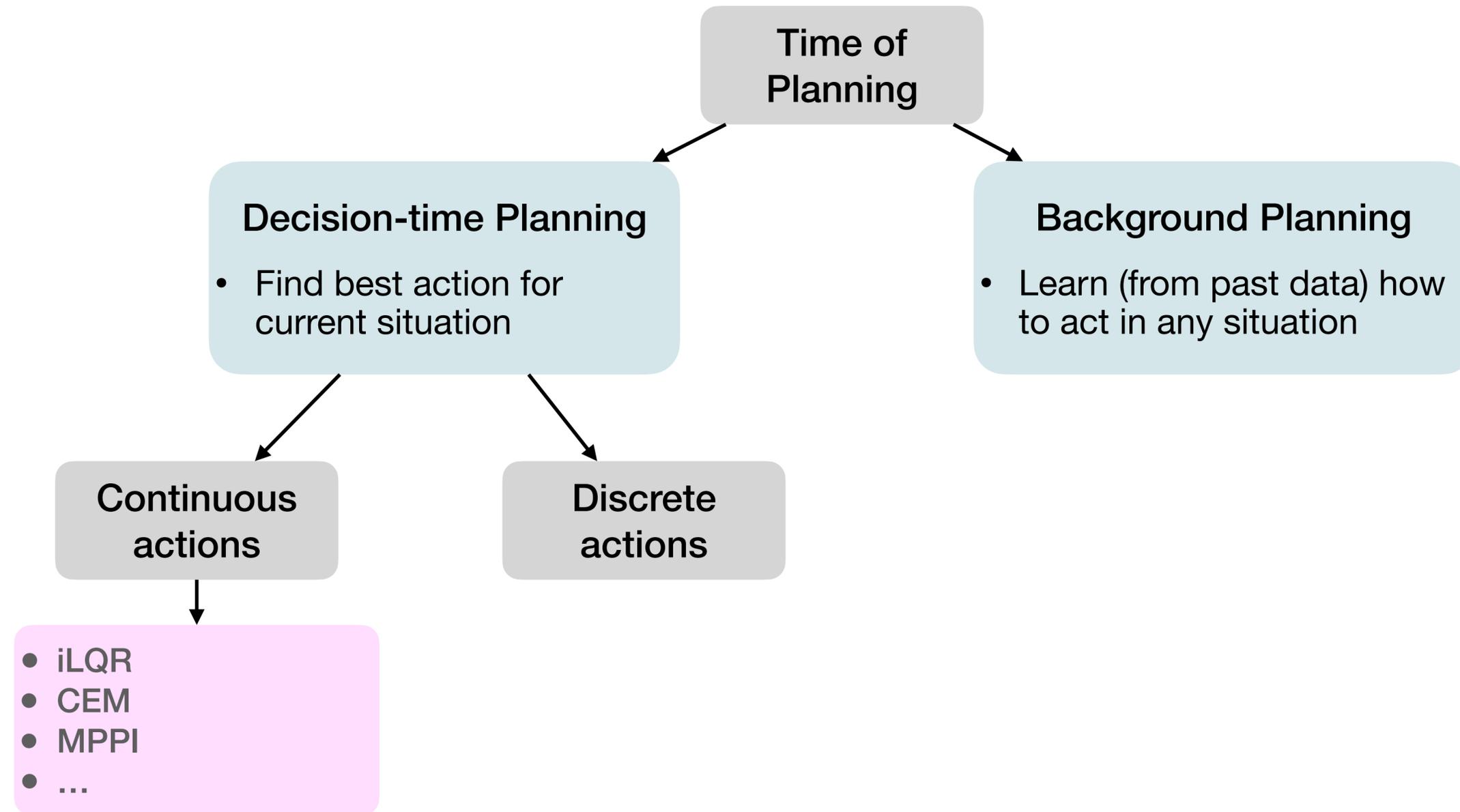
# How Do We Use The "Model"?

## Background vs Decision-time Planning



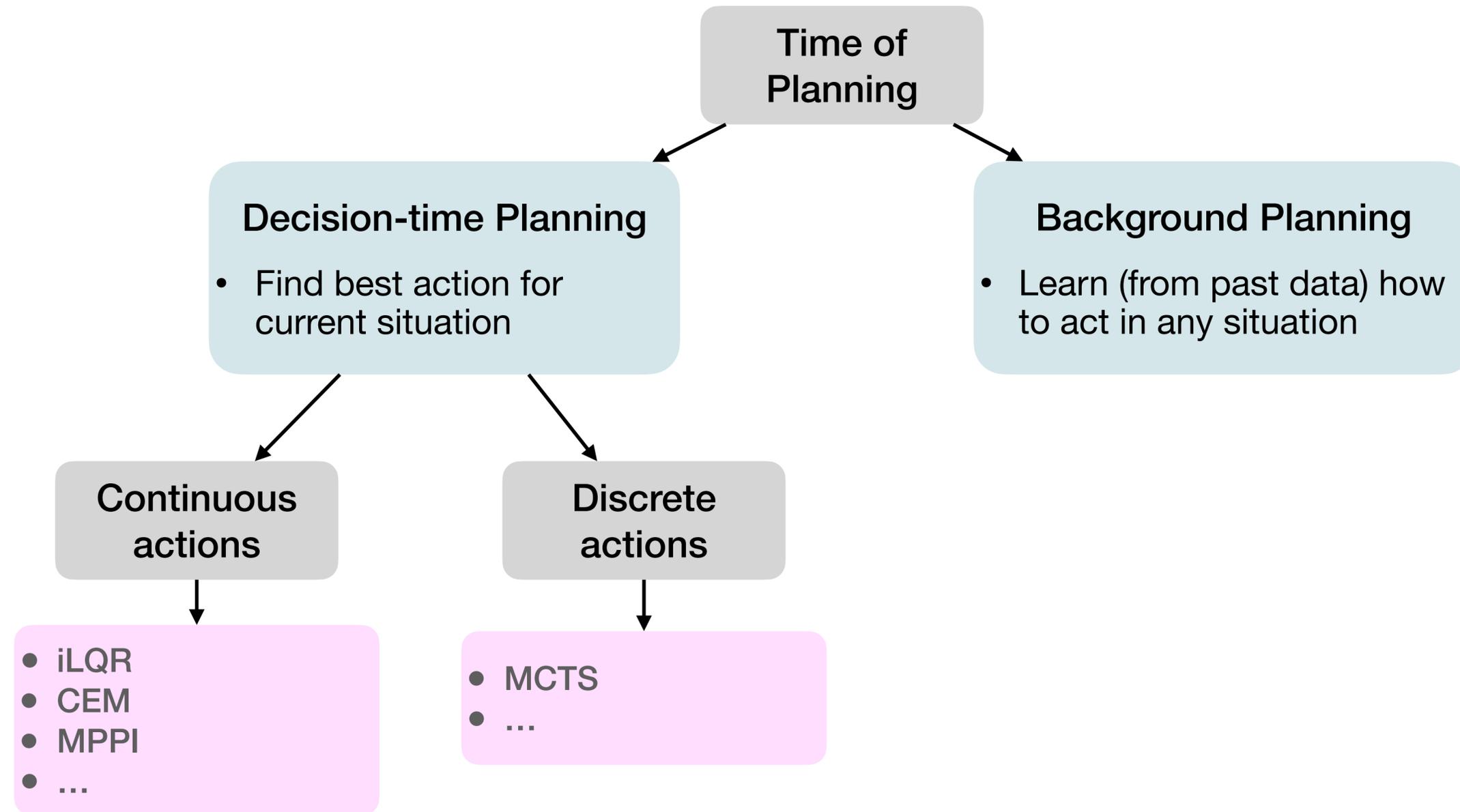
# How Do We Use The "Model"?

## Background vs Decision-time Planning



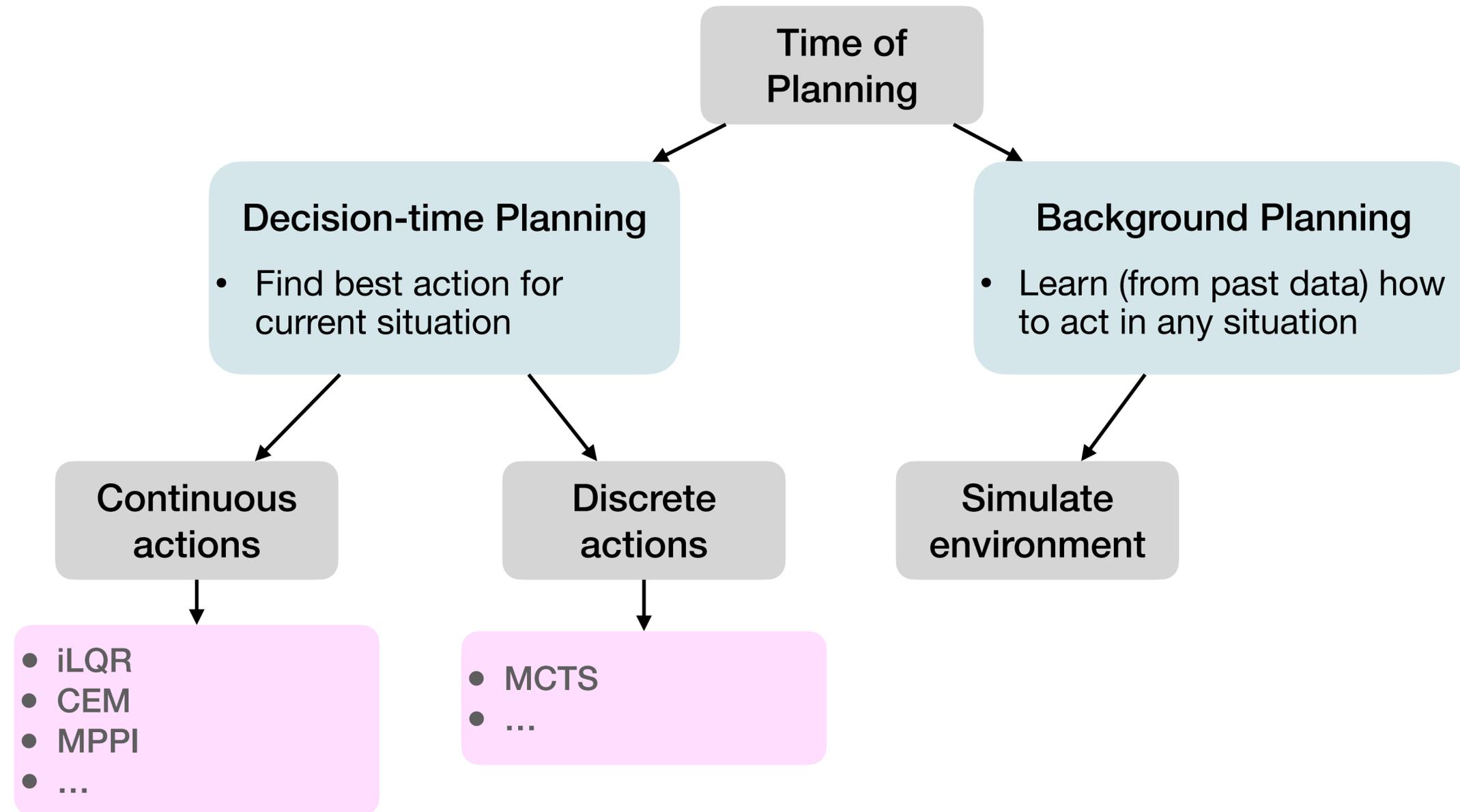
# How Do We Use The "Model"?

## Background vs Decision-time Planning



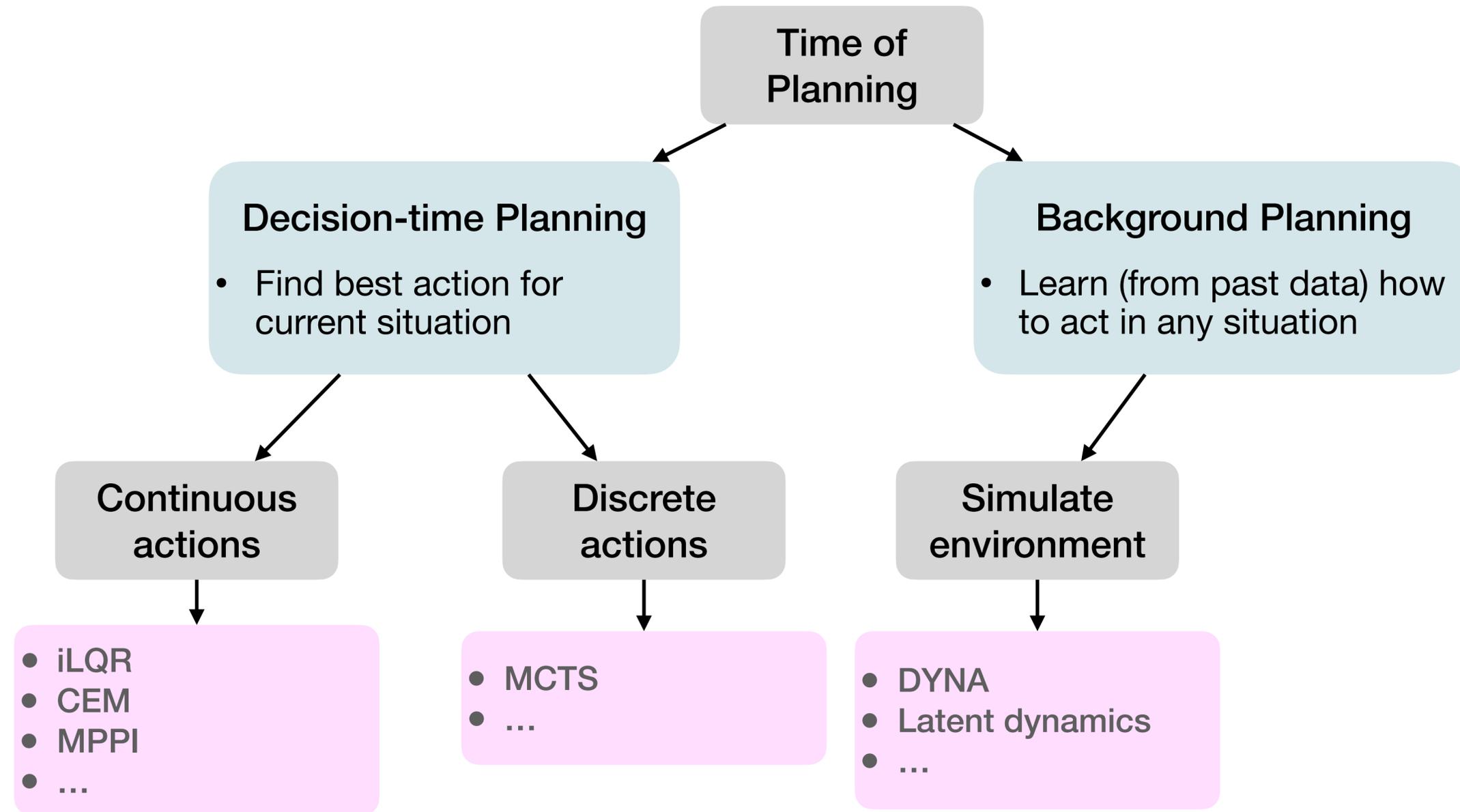
# How Do We Use The "Model"?

## Background vs Decision-time Planning



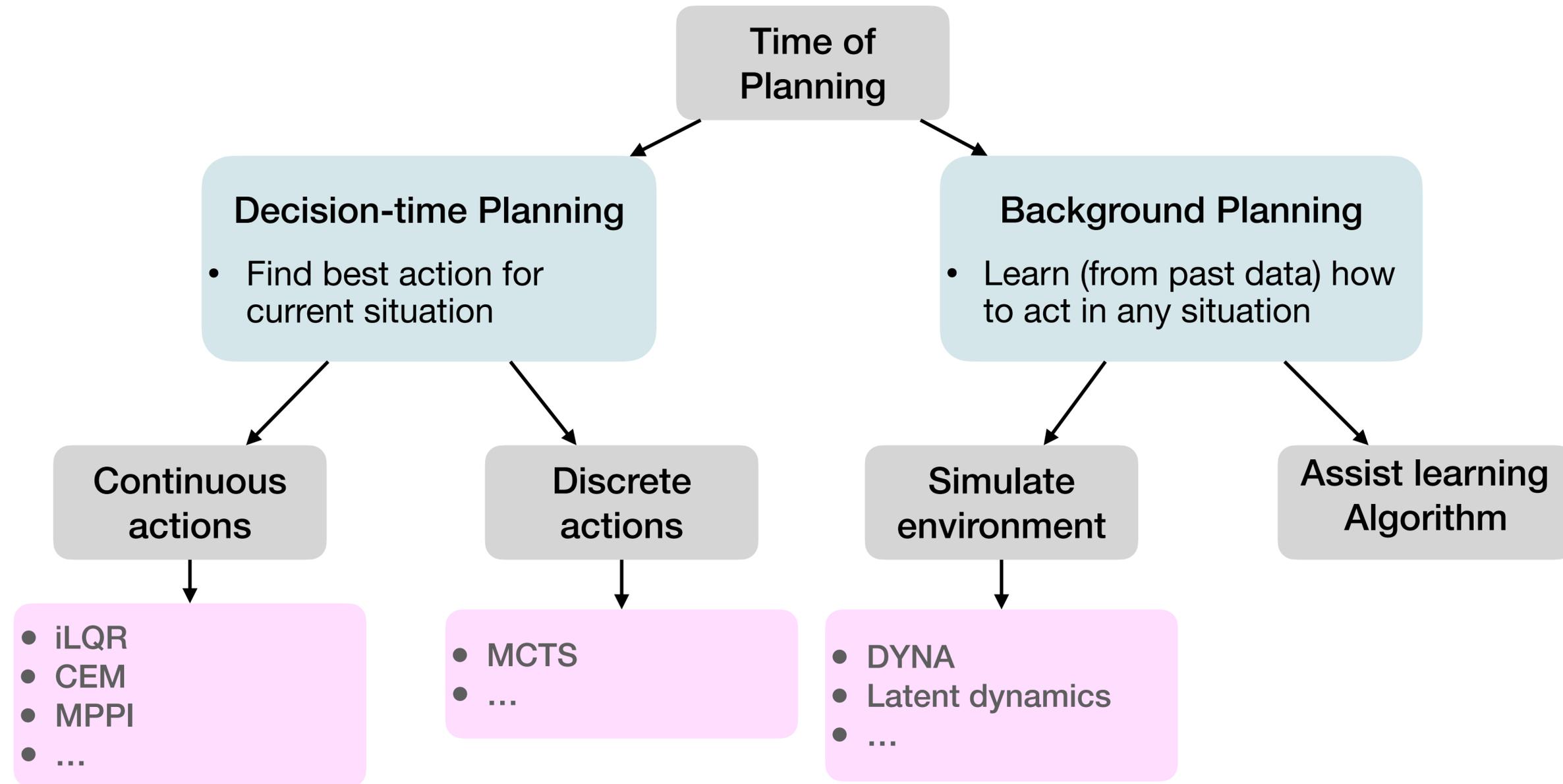
# How Do We Use The "Model"?

## Background vs Decision-time Planning



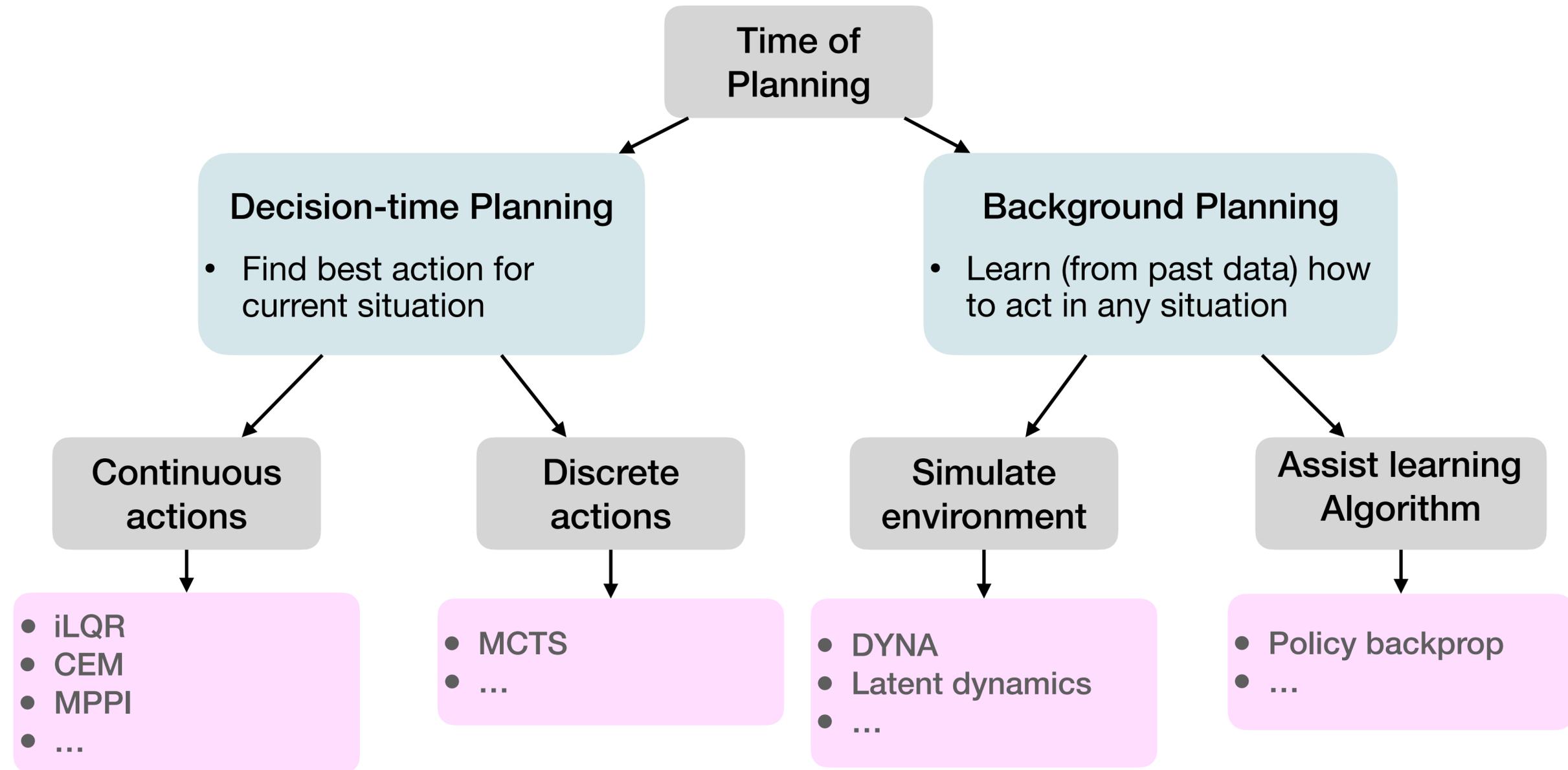
# How Do We Use The "Model"?

## Background vs Decision-time Planning



# How Do We Use The "Model"?

## Background vs Decision-time Planning



# Background vs Decision-time Planning

**Background planning**

**FCAI**

**Decision-time planning**

**fcai.fi**

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*

**FCAI**

## Decision-time planning

**fcai.fi**

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*

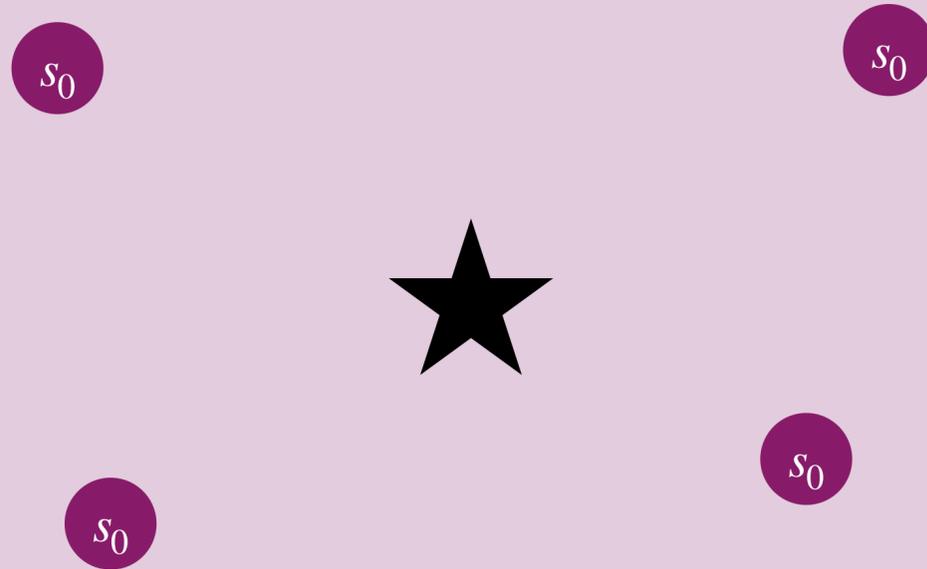


## Decision-time planning

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*

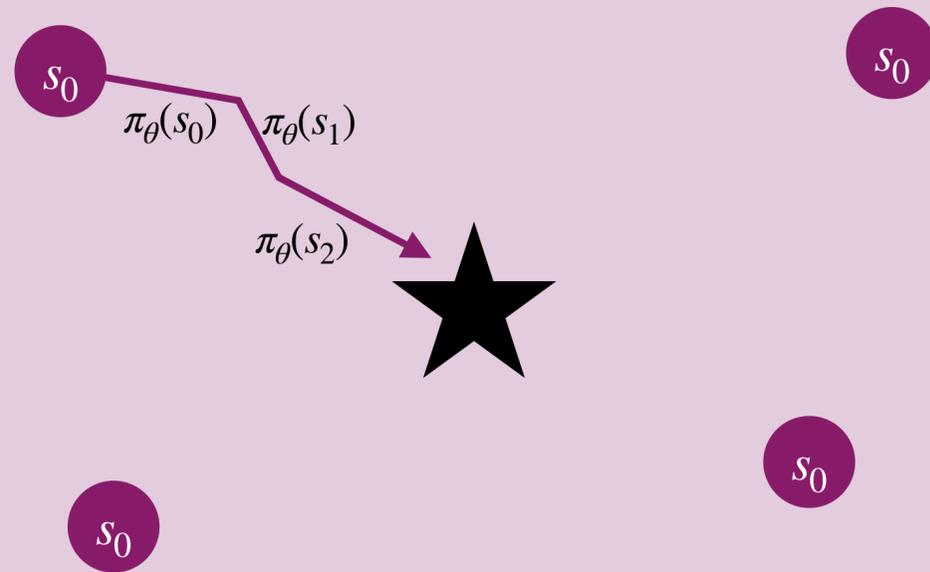


## Decision-time planning

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*

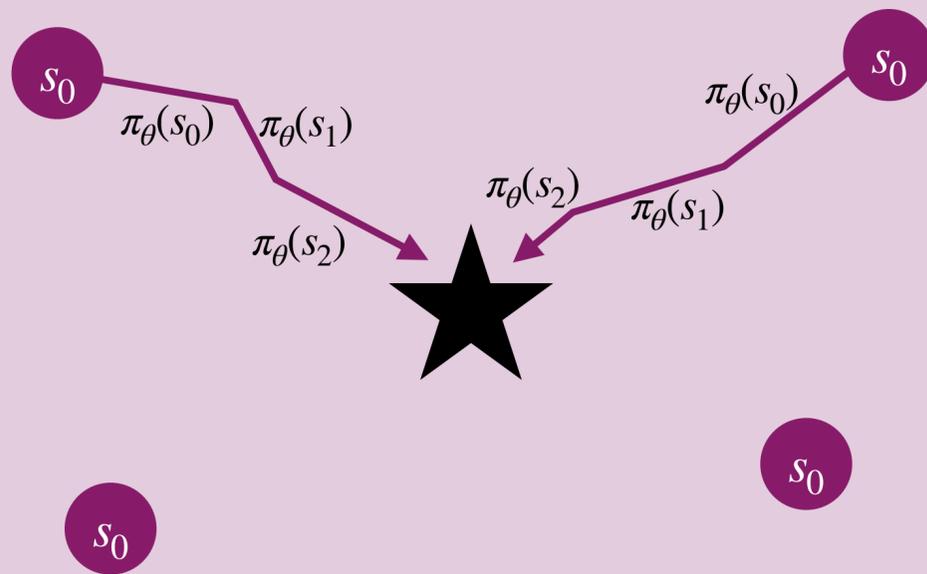


## Decision-time planning

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*

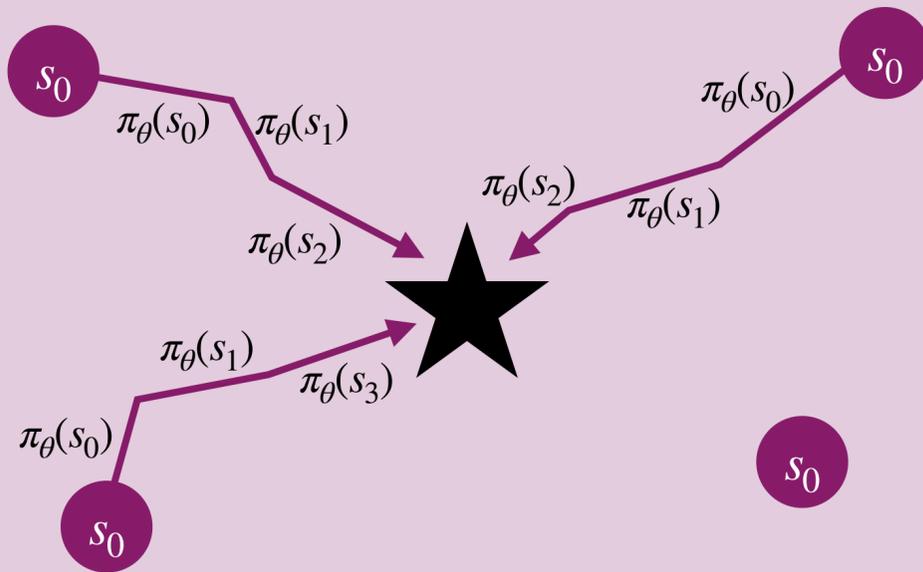


## Decision-time planning

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*

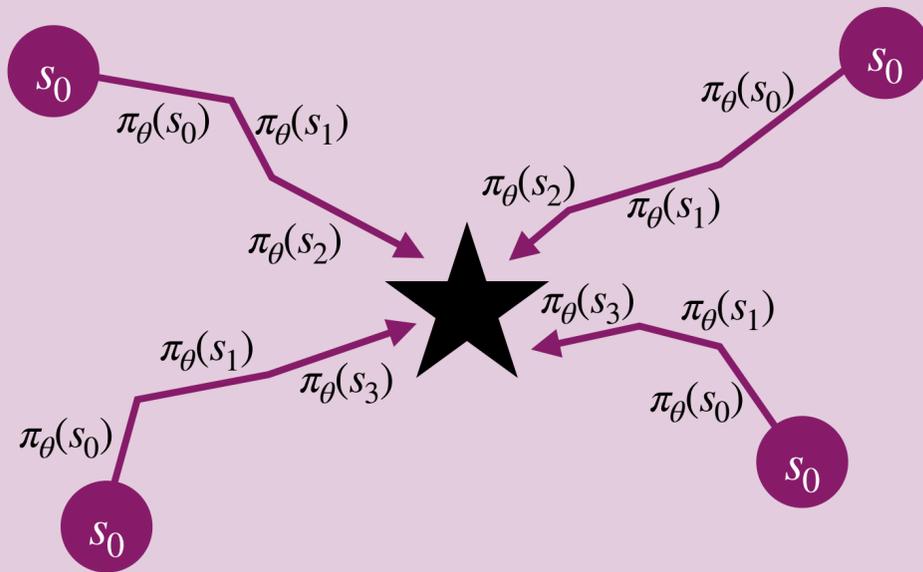


## Decision-time planning

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*

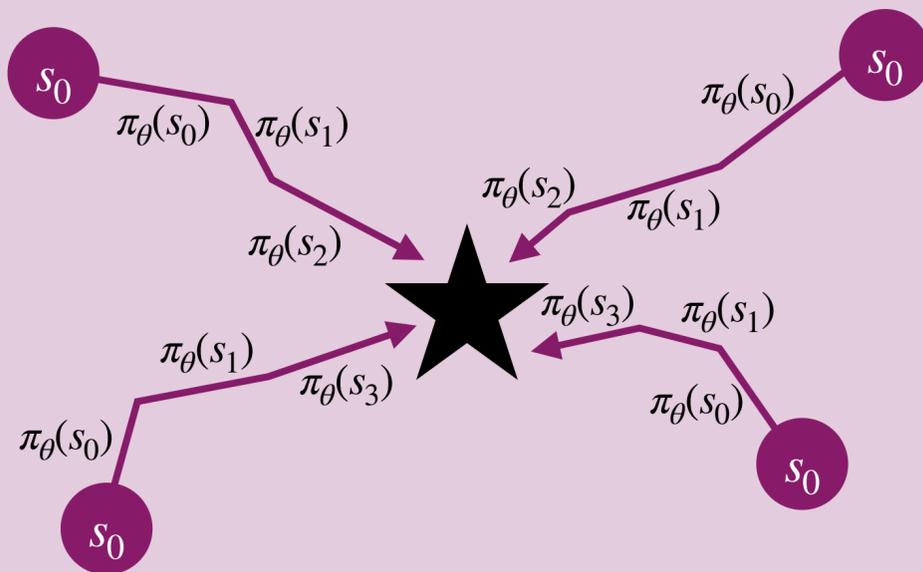


## Decision-time planning

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



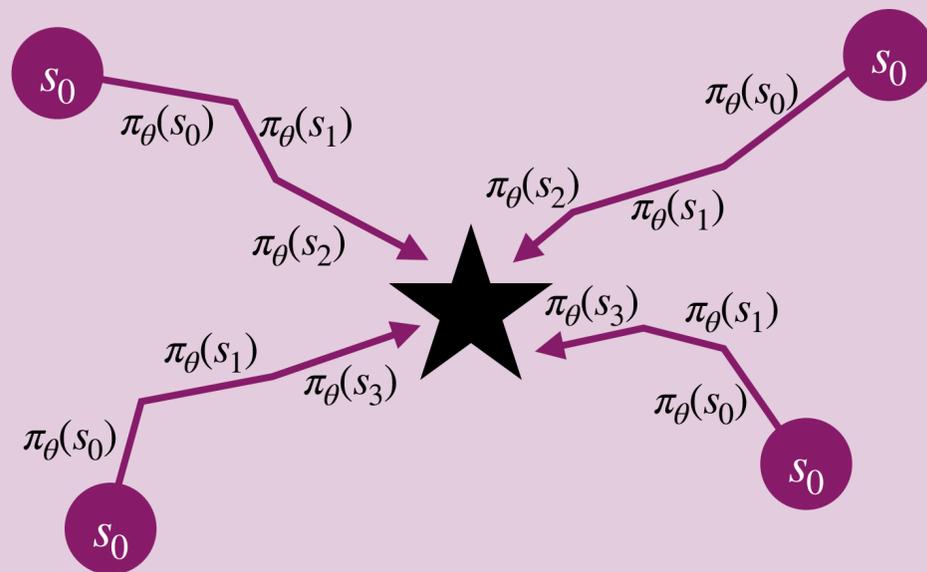
Optimisation variables:  $\theta$

## Decision-time planning

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



**Optimisation variables:  $\theta$**

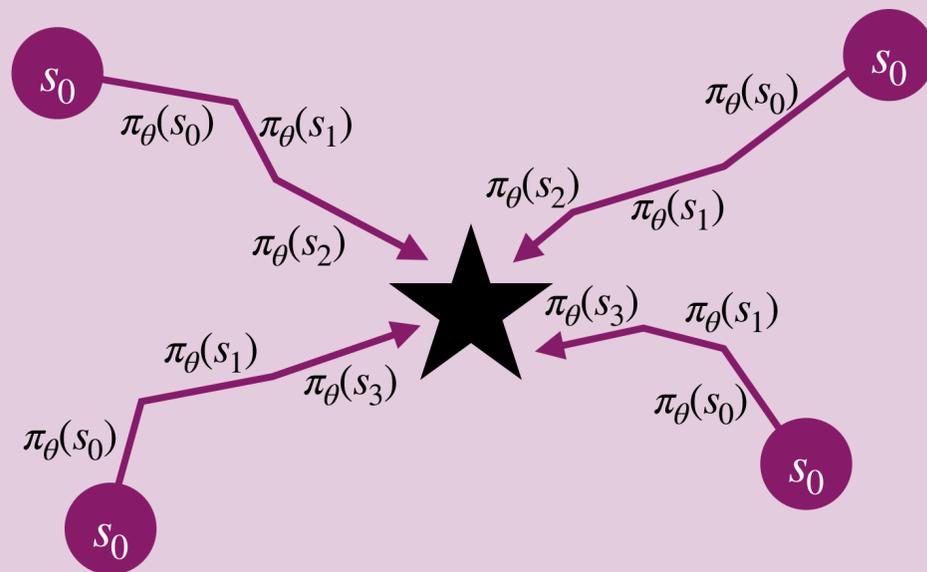
Parameters of policy  $\pi_\theta(s)$ , value  $Q_\theta(s, a)$ , etc

## Decision-time planning

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



**Optimisation variables:  $\theta$**

Parameters of policy  $\pi_\theta(s)$ , value  $Q_\theta(s, a)$ , etc

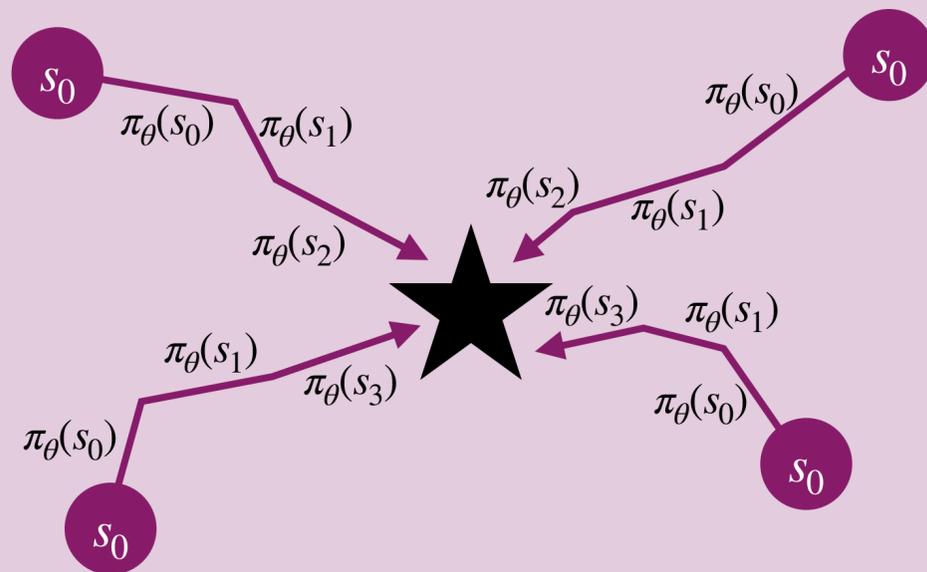
$$J(\theta) = \mathbb{E}_{s_0} \left[ \sum_{t=0}^H r(s_t, \pi_\theta(s_t)) \right]$$

## Decision-time planning

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



**Optimisation variables:  $\theta$**

Parameters of policy  $\pi_\theta(s)$ , value  $Q_\theta(s, a)$ , etc

$$J(\theta) = \mathbb{E}_{s_0} \left[ \sum_{t=0}^H r(s_t, \pi_\theta(s_t)) \right]$$

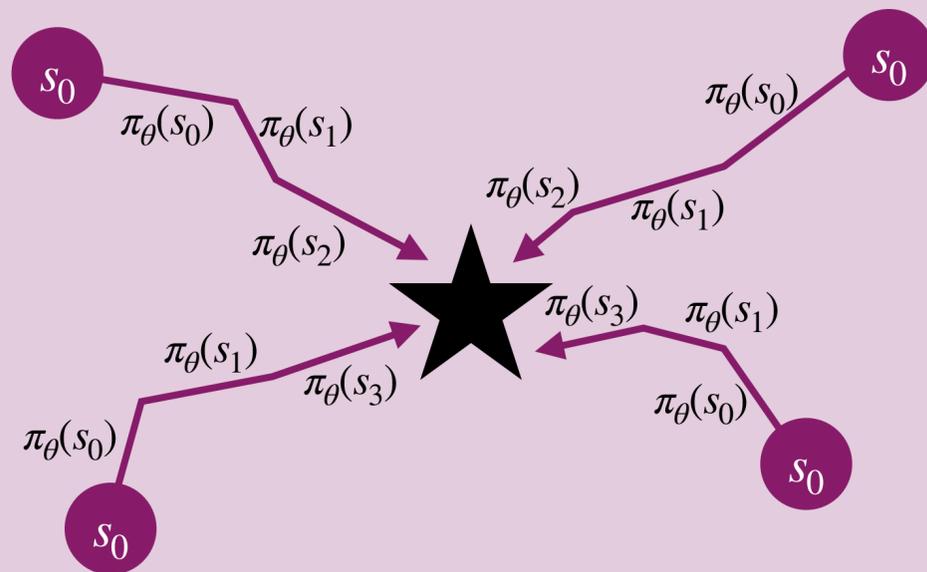
## Decision-time planning

*Find best action for current situation*

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



**Optimisation variables:  $\theta$**

Parameters of policy  $\pi_\theta(s)$ , value  $Q_\theta(s, a)$ , etc

$$J(\theta) = \mathbb{E}_{s_0} \left[ \sum_{t=0}^H r(s_t, \pi_\theta(s_t)) \right]$$

## Decision-time planning

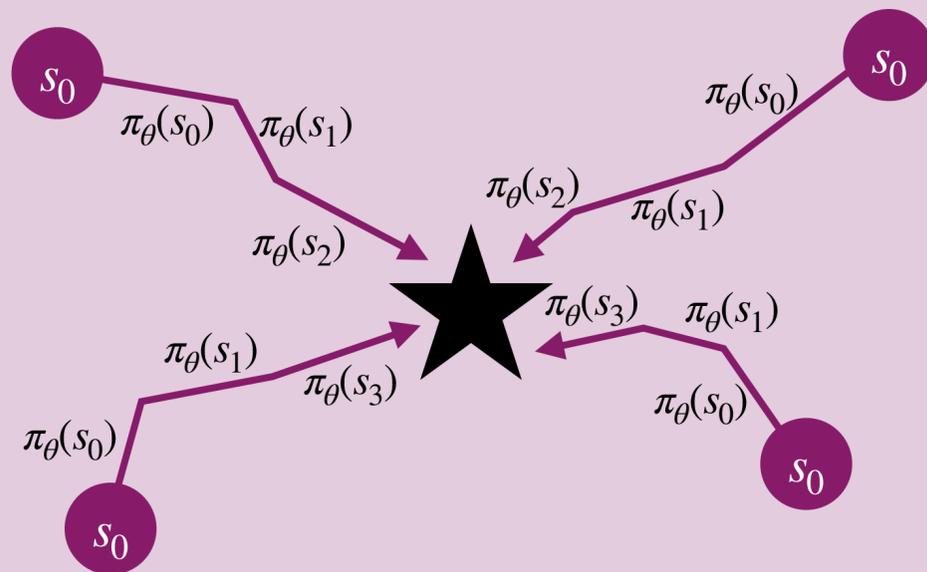
*Find best action for current situation*



# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



**Optimisation variables:  $\theta$**

Parameters of policy  $\pi_\theta(s)$ , value  $Q_\theta(s, a)$ , etc

$$J(\theta) = \mathbb{E}_{s_0} \left[ \sum_{t=0}^H r(s_t, \pi_\theta(s_t)) \right]$$

## Decision-time planning

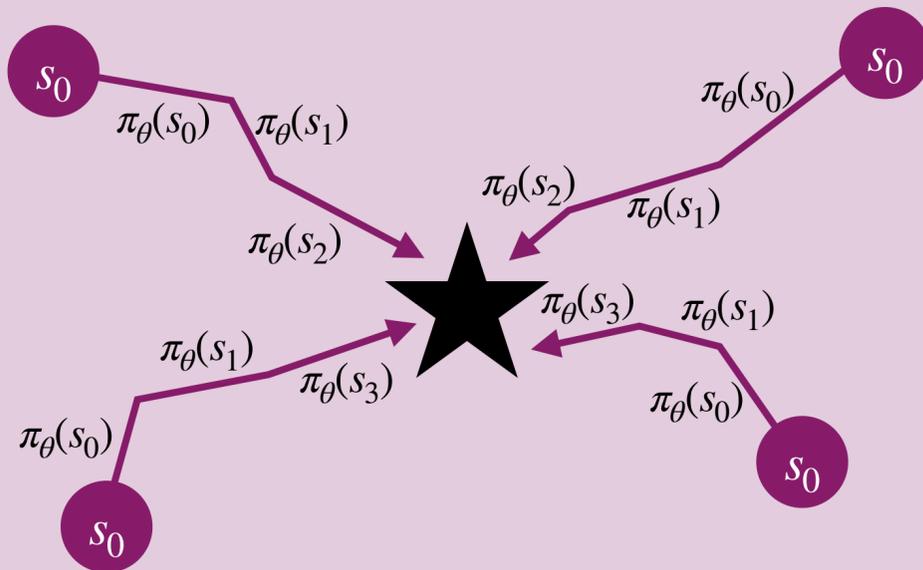
*Find best action for current situation*



# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



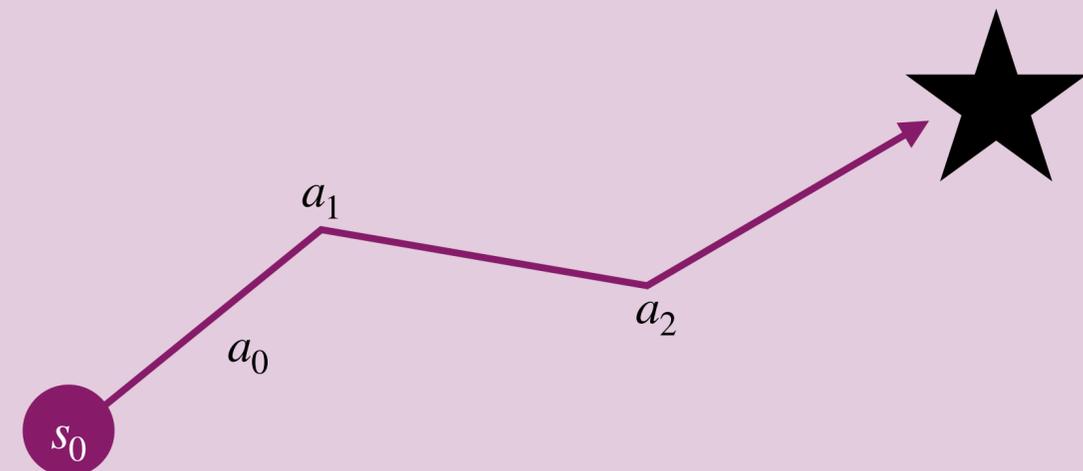
**Optimisation variables:  $\theta$**

Parameters of policy  $\pi_\theta(s)$ , value  $Q_\theta(s, a)$ , etc

$$J(\theta) = \mathbb{E}_{s_0} \left[ \sum_{t=0}^H r(s_t, \pi_\theta(s_t)) \right]$$

## Decision-time planning

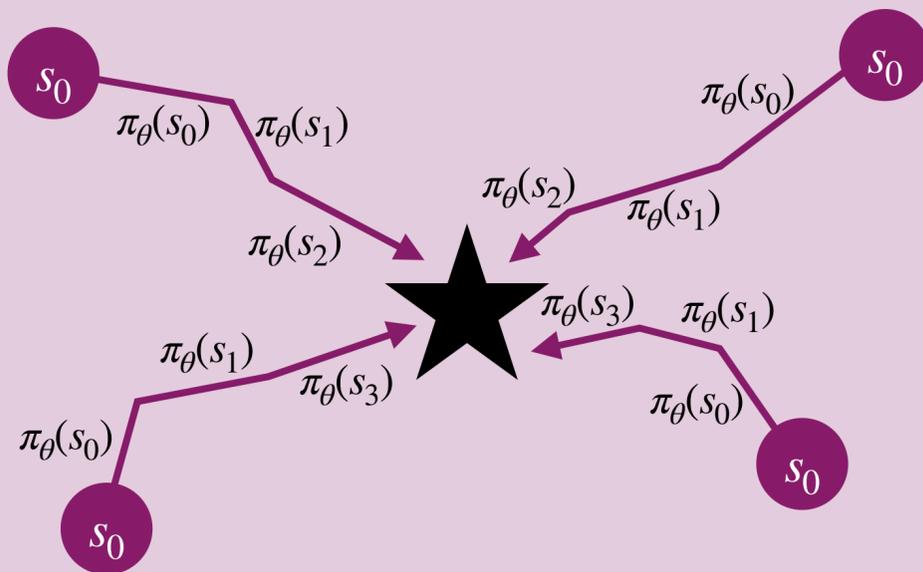
*Find best action for current situation*



# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



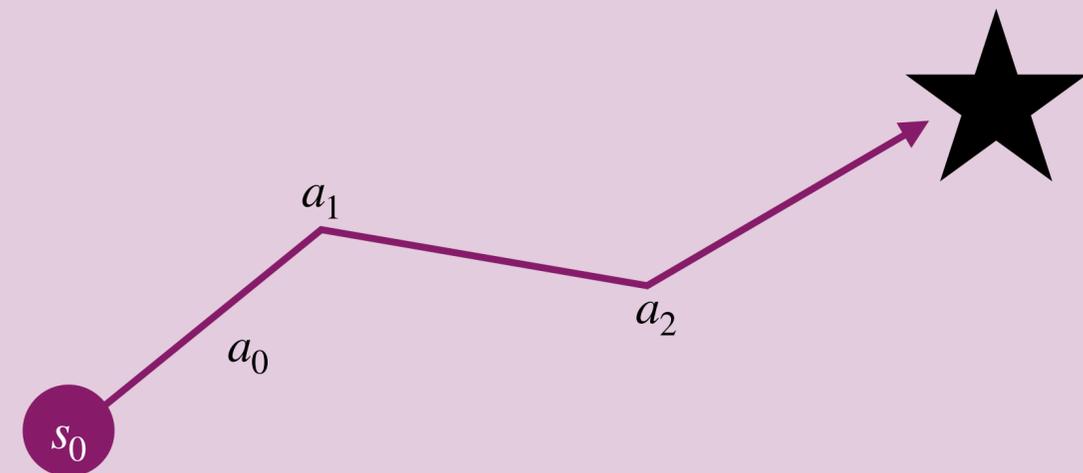
**Optimisation variables:**  $\theta$

Parameters of policy  $\pi_\theta(s)$ , value  $Q_\theta(s, a)$ , etc

$$J(\theta) = \mathbb{E}_{s_0} \left[ \sum_{t=0}^H r(s_t, \pi_\theta(s_t)) \right]$$

## Decision-time planning

*Find best action for current situation*

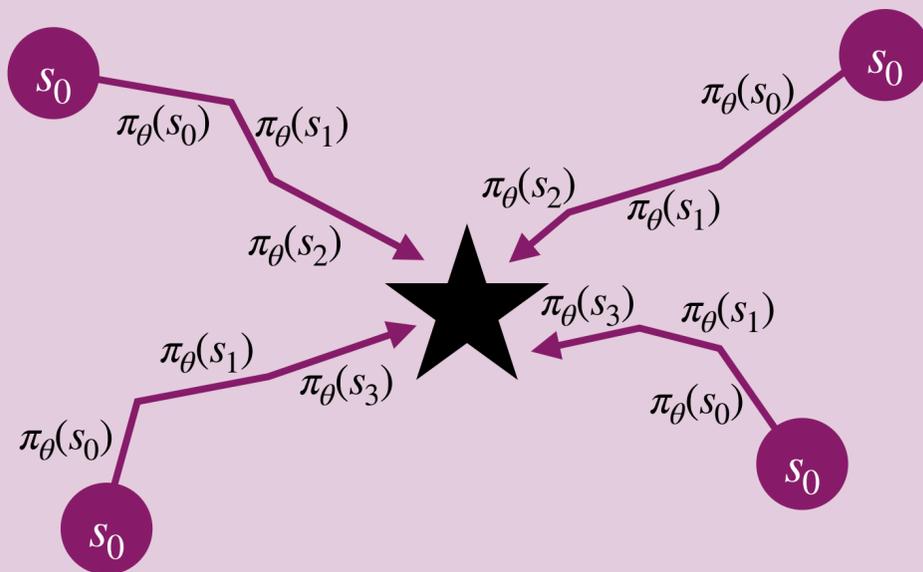


**Optimisation variables:**  $a_0, \dots, a_H$

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



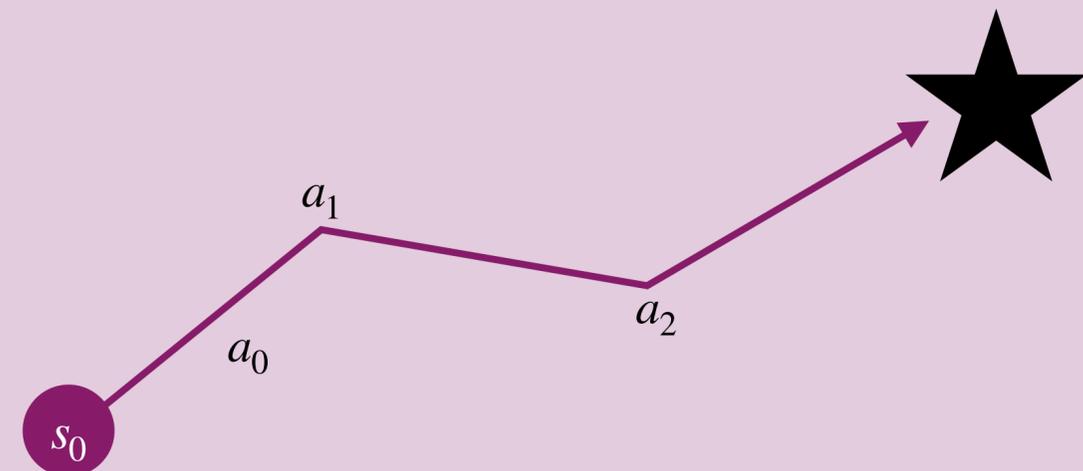
**Optimisation variables:**  $\theta$

Parameters of policy  $\pi_\theta(s)$ , value  $Q_\theta(s, a)$ , etc

$$J(\theta) = \mathbb{E}_{s_0} \left[ \sum_{t=0}^H r(s_t, \pi_\theta(s_t)) \right]$$

## Decision-time planning

*Find best action for current situation*



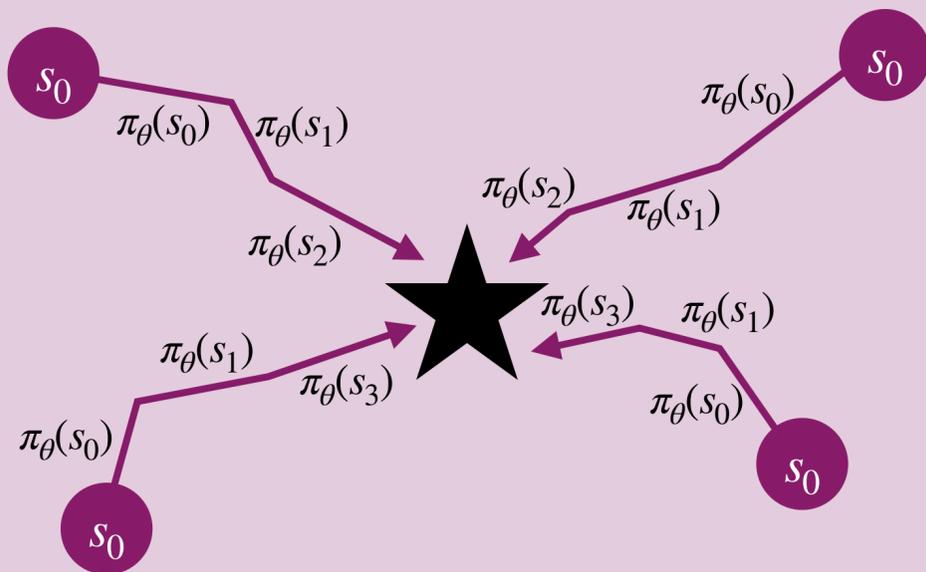
**Optimisation variables:**  $a_0, \dots, a_H$

Sequence of actions (and maybe also states)

# Background vs Decision-time Planning

## Background planning

*Learn how to act in any situation*



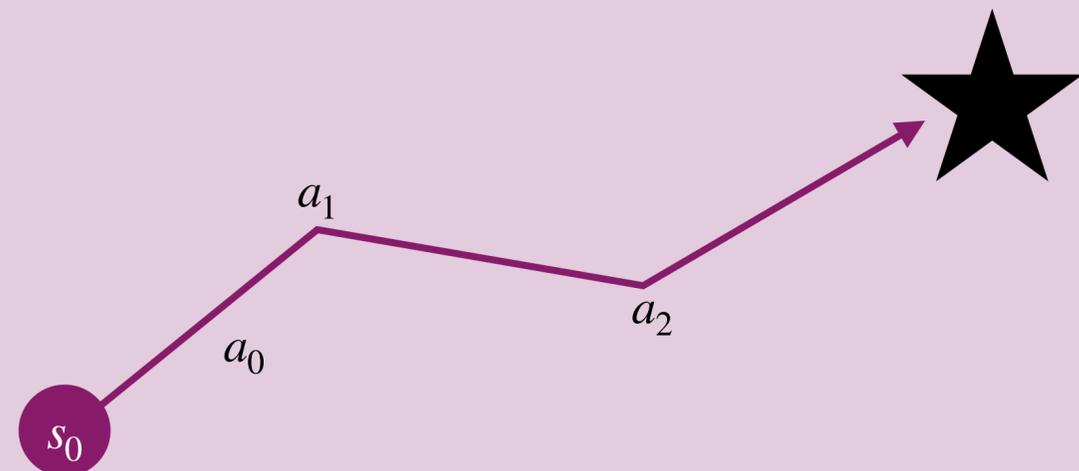
**Optimisation variables:**  $\theta$

Parameters of policy  $\pi_\theta(s)$ , value  $Q_\theta(s, a)$ , etc

$$J(\theta) = \mathbb{E}_{s_0} \left[ \sum_{t=0}^H r(s_t, \pi_\theta(s_t)) \right]$$

## Decision-time planning

*Find best action for current situation*



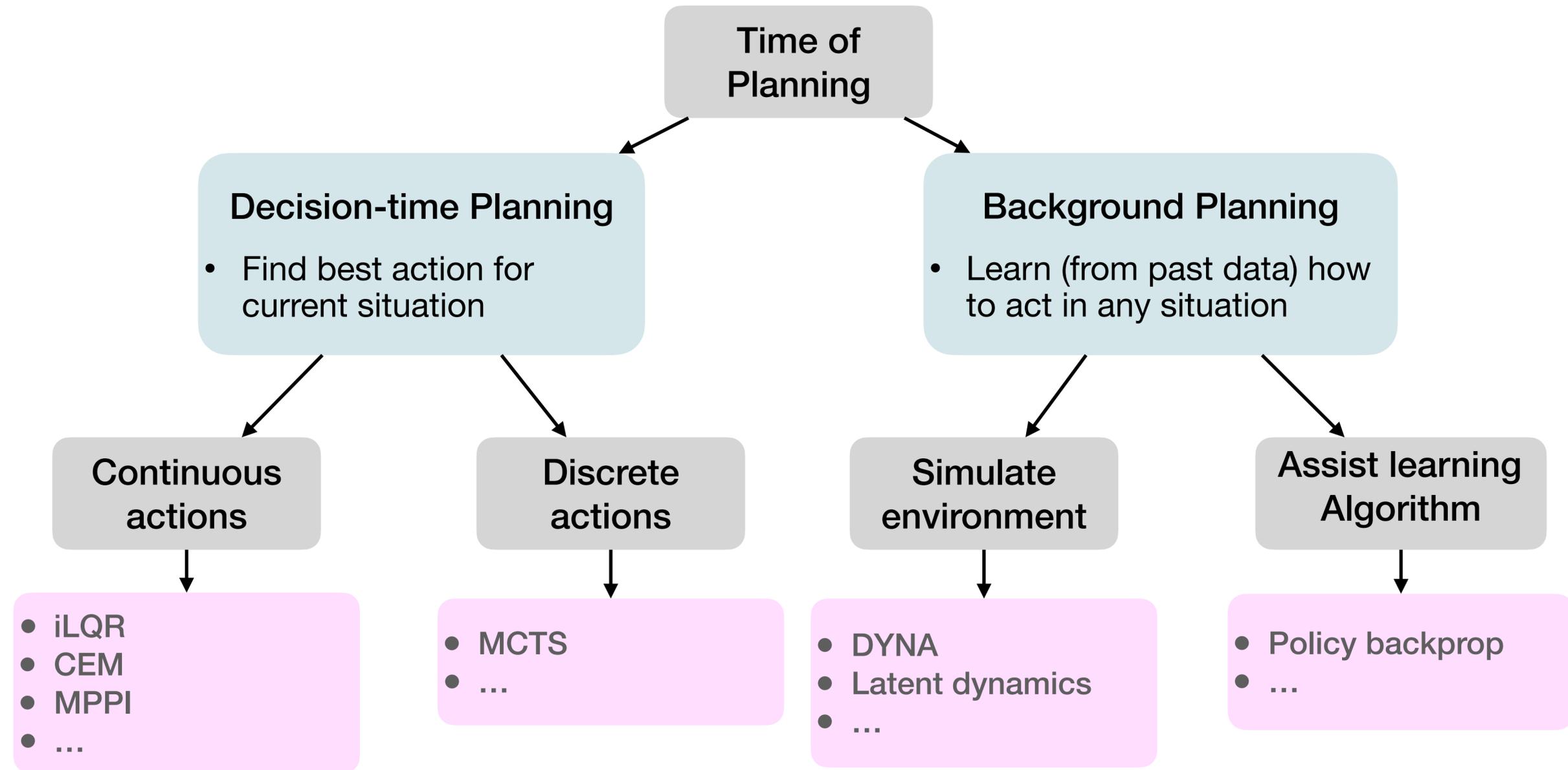
**Optimisation variables:**  $a_0, \dots, a_H$

Sequence of actions (and maybe also states)

$$J(a_0, \dots, a_H) = \sum_{t=0}^H r(s_t, a_t)$$

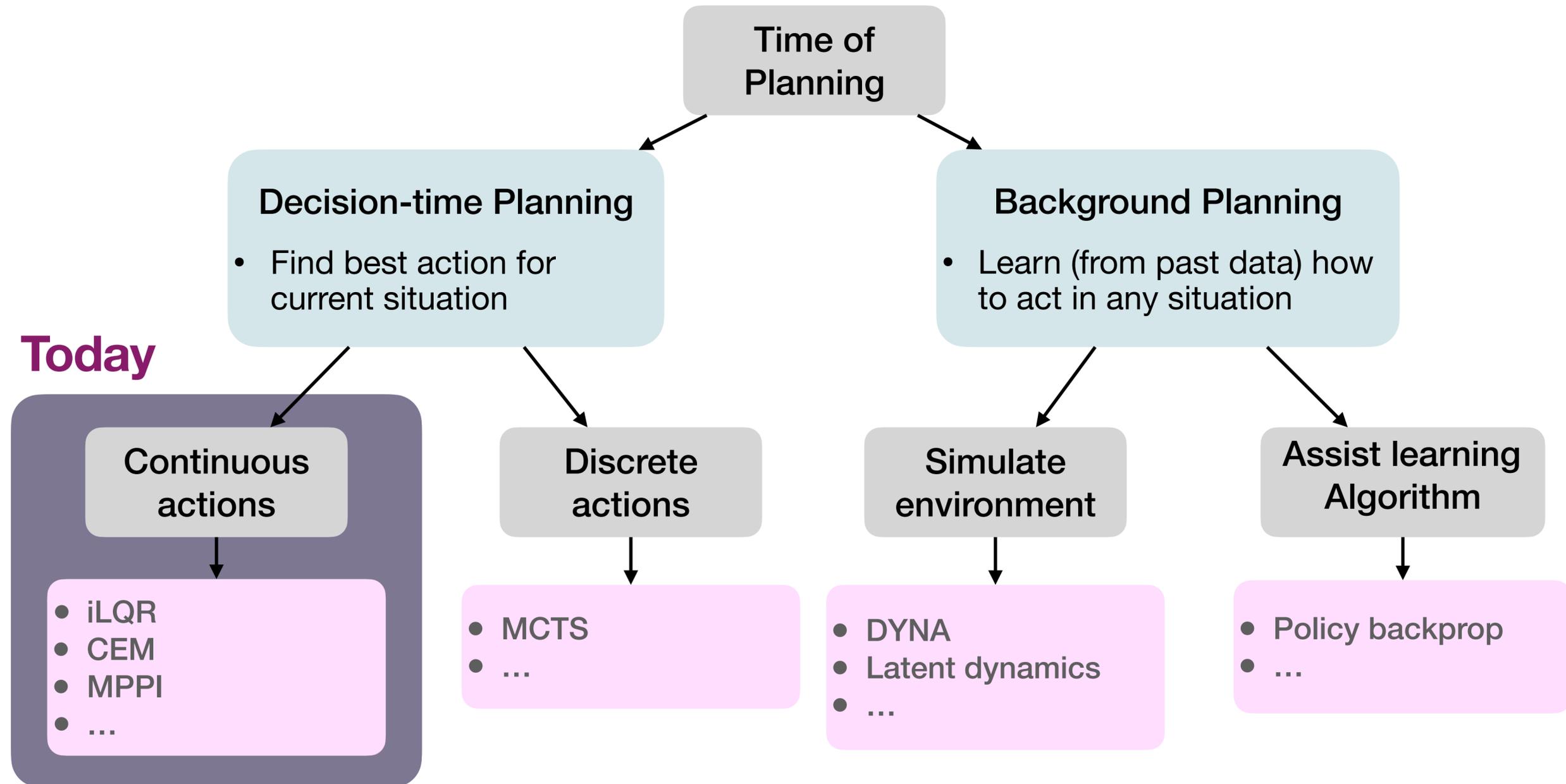
# How Do We Use The "Model"?

## Background vs Decision-time Planning



# How Do We Use The "Model"?

## Background vs Decision-time Planning



**Today**

# Decision-time Planning (Continuous Actions)

# Decision-time Planning (Continuous Actions)

We'll start by assuming known, deterministic dynamics

# Decision-time Planning (Continuous Actions)

We'll start by assuming known, deterministic dynamics

$$s_{t+1} = f(s_t, a_t)$$

# Decision-time Planning

## Trajectory optimisation

# Decision-time Planning

## Trajectory optimisation

Observe state  $s$

# Decision-time Planning

## Trajectory optimisation

Observe state  $s$

Plan  $a_0, \dots, a_H$  to maximise return  $\sum_{t=0}^H \gamma^t r(s_t, a_t)$  s.t.  $s_0 = s$

# Decision-time Planning

## Trajectory optimisation

Observe state  $s$

Plan  $a_0, \dots, a_H$  to maximise return  $\sum_{t=0}^H \gamma^t r(s_t, a_t)$  s.t.  $s_0 = s$

Execute each action

# Trajectory Optimisation

**Shooting methods**

**FCAI**

**Collocation methods**

**fcai.fi**

# Trajectory Optimisation

## Shooting methods

Optimisation variables:  $a_0, \dots, a_H$

Actions

$$J(a_{0:H}) = \sum_{t=0}^H r(s_t, a_t)$$

## Collocation methods

# Trajectory Optimisation

## Shooting methods

Optimisation variables:  $a_0, \dots, a_H$

Actions

$$J(a_{0:H}) = \sum_{t=0}^H r(s_t, a_t)$$

## Collocation methods

Optimisation variables:  $a_0, s_0, \dots, a_H, s_H$

Actions and states

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H r(s_t, a_t)$$

$$\text{s.t. } \|s_{t+1} - f(s_t, a_t)\| = 0$$

# Shooting Methods

## Illustration



# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$



# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions

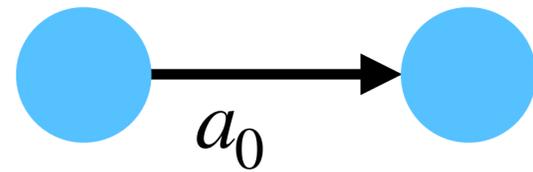


# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions

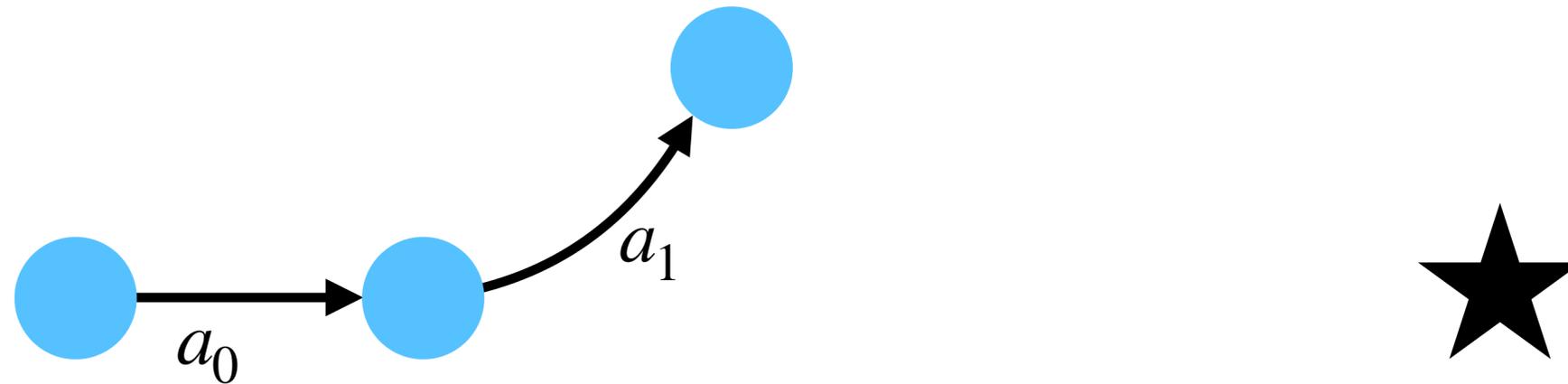


# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions

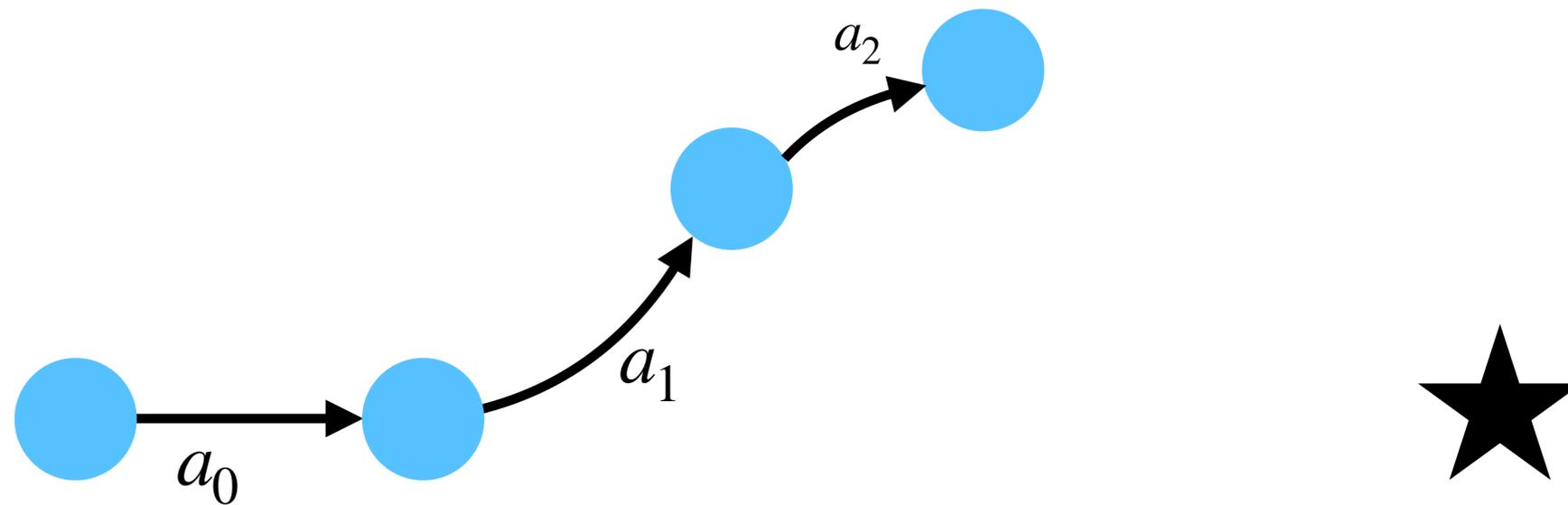


# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions

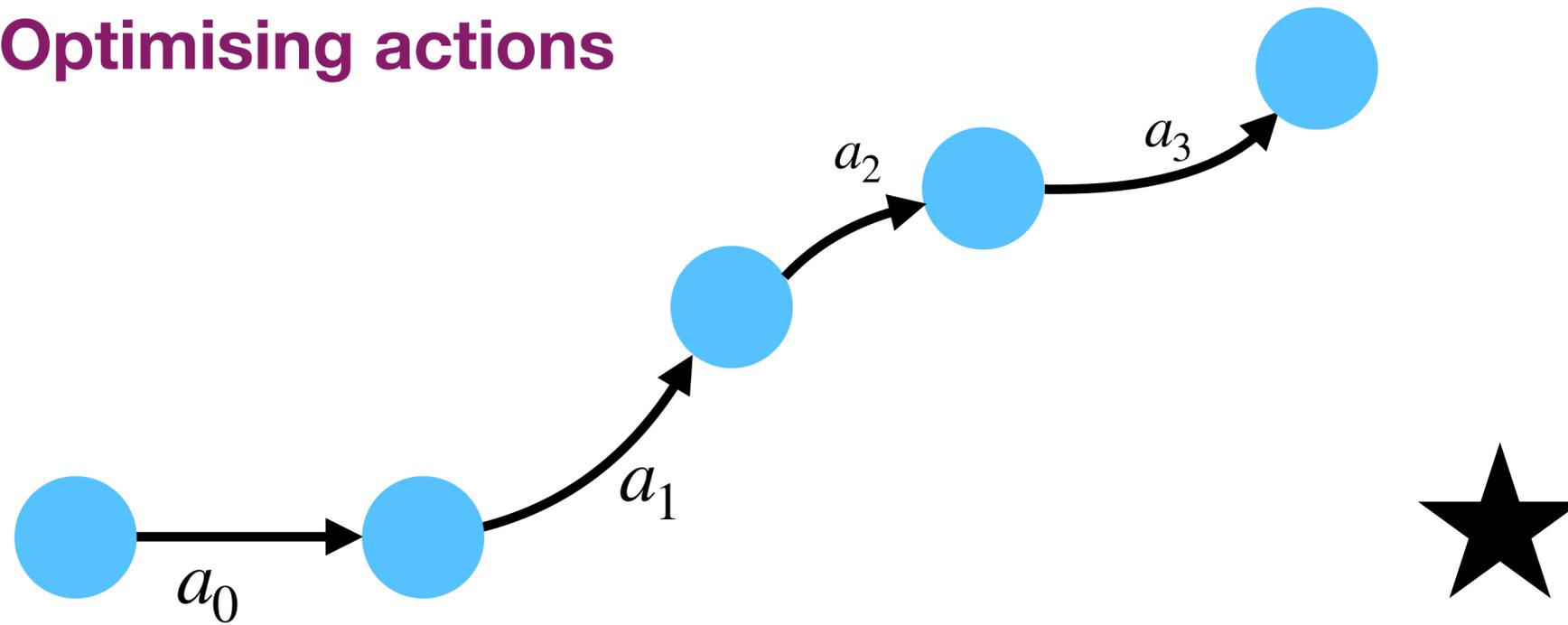


# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions

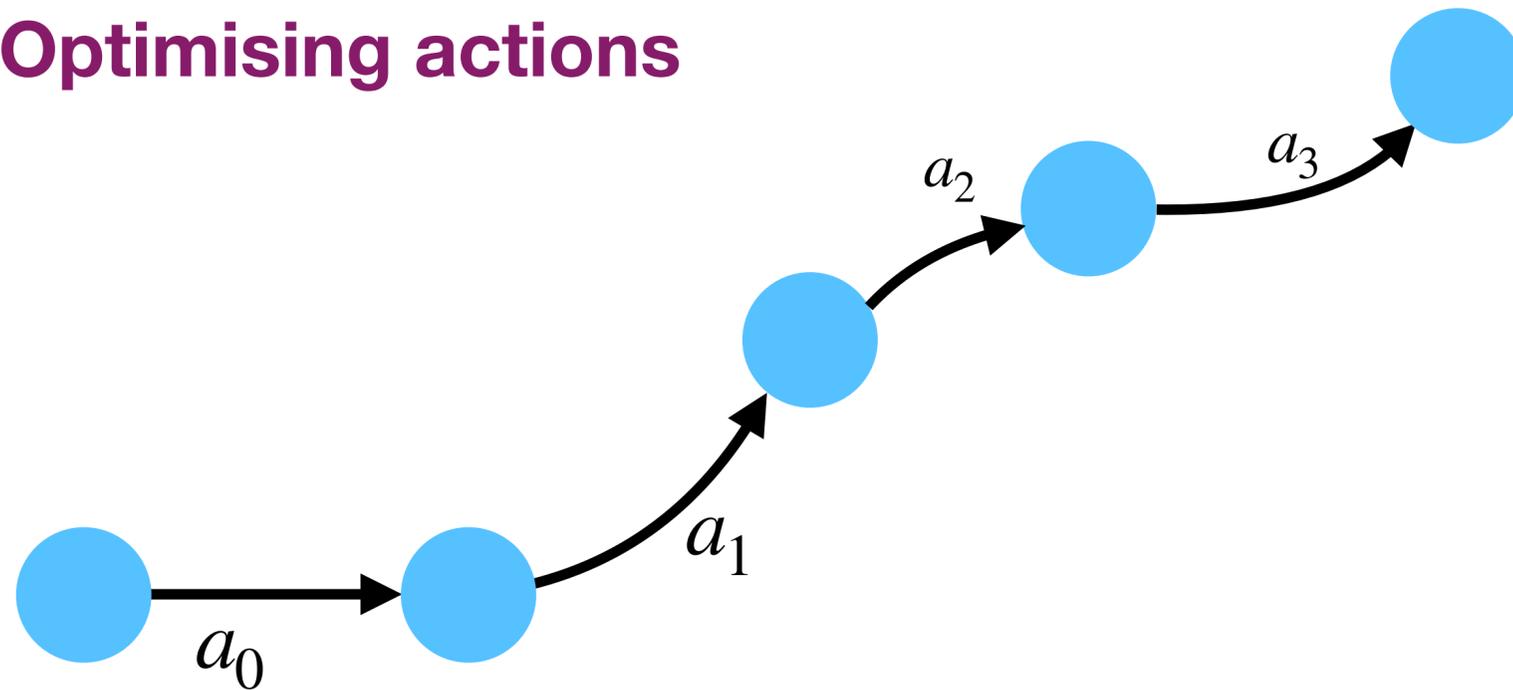


# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



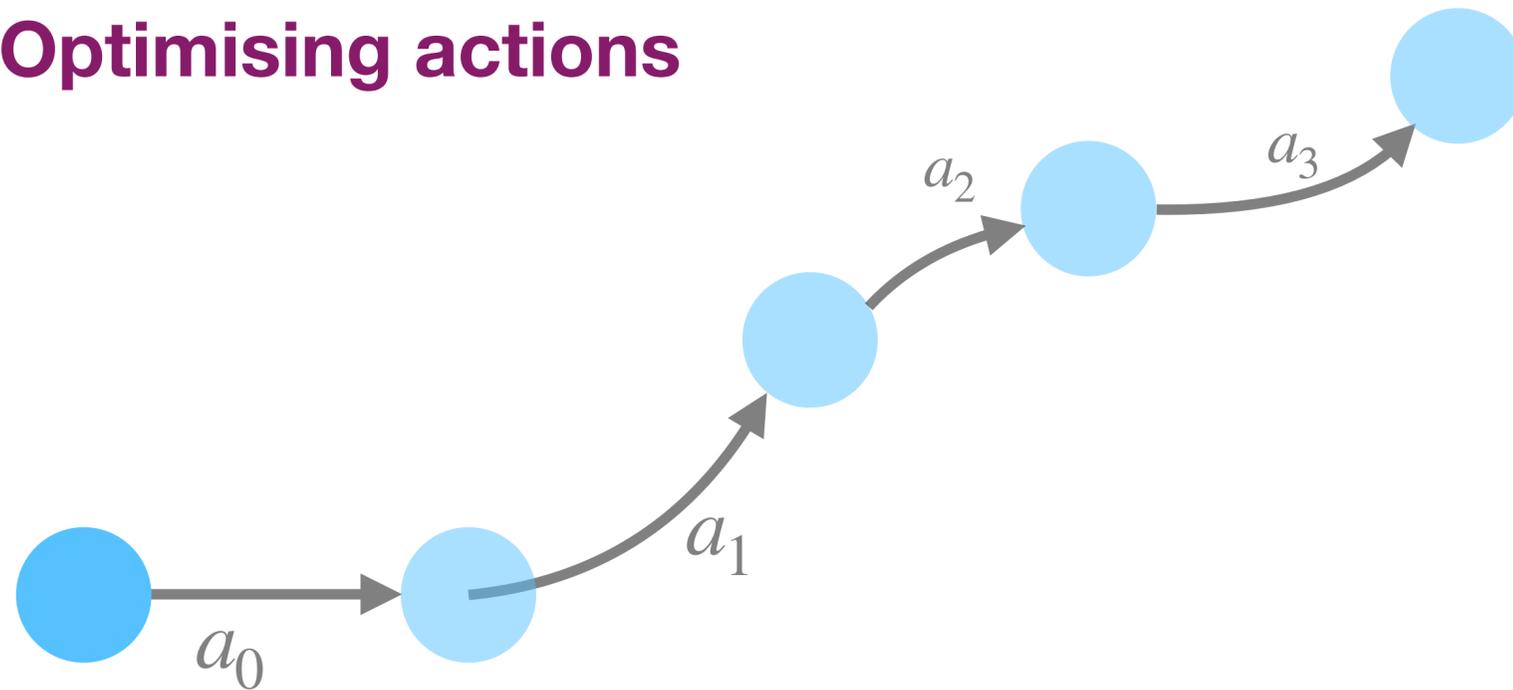
Recursively evaluate dynamics

# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



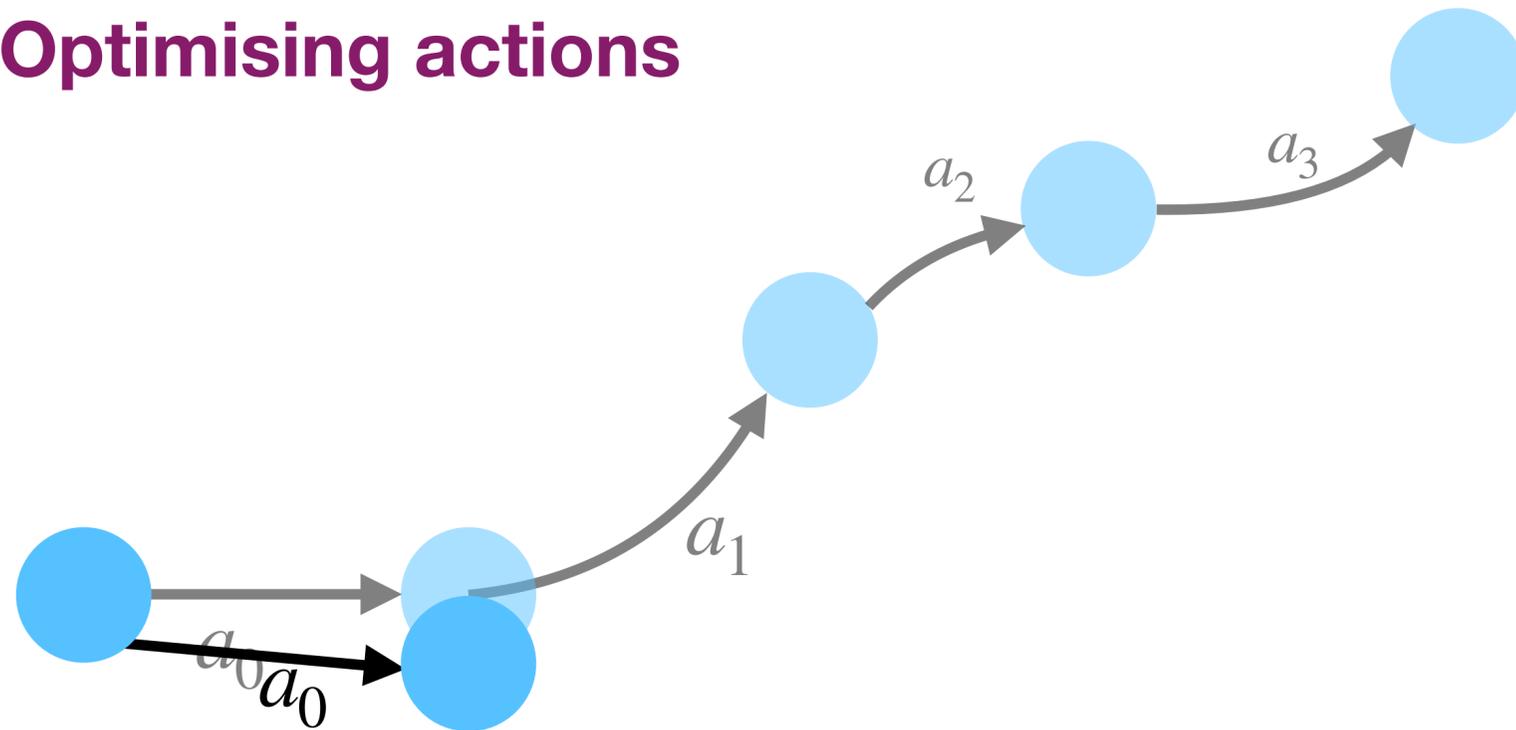
Recursively evaluate dynamics

# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



Recursively evaluate dynamics

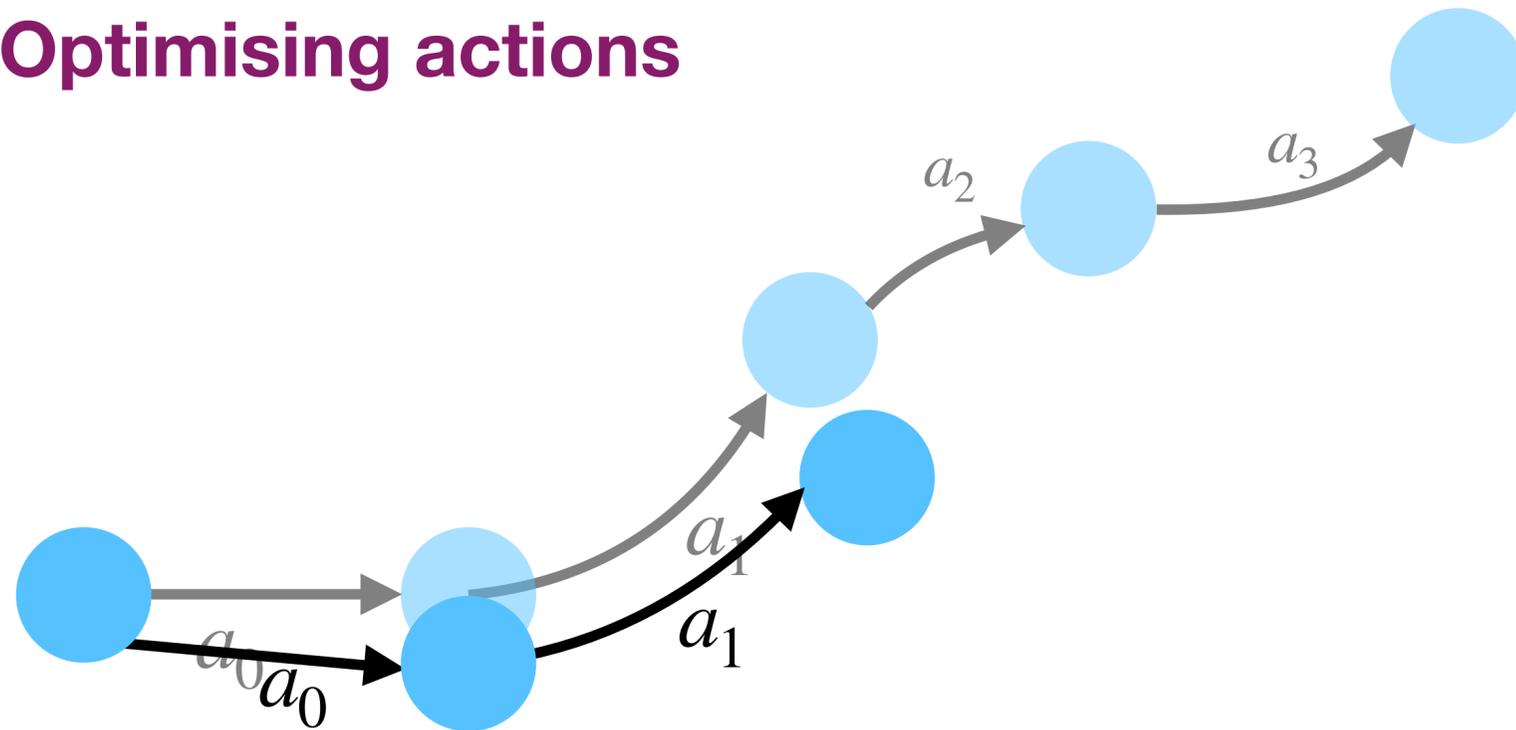


# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



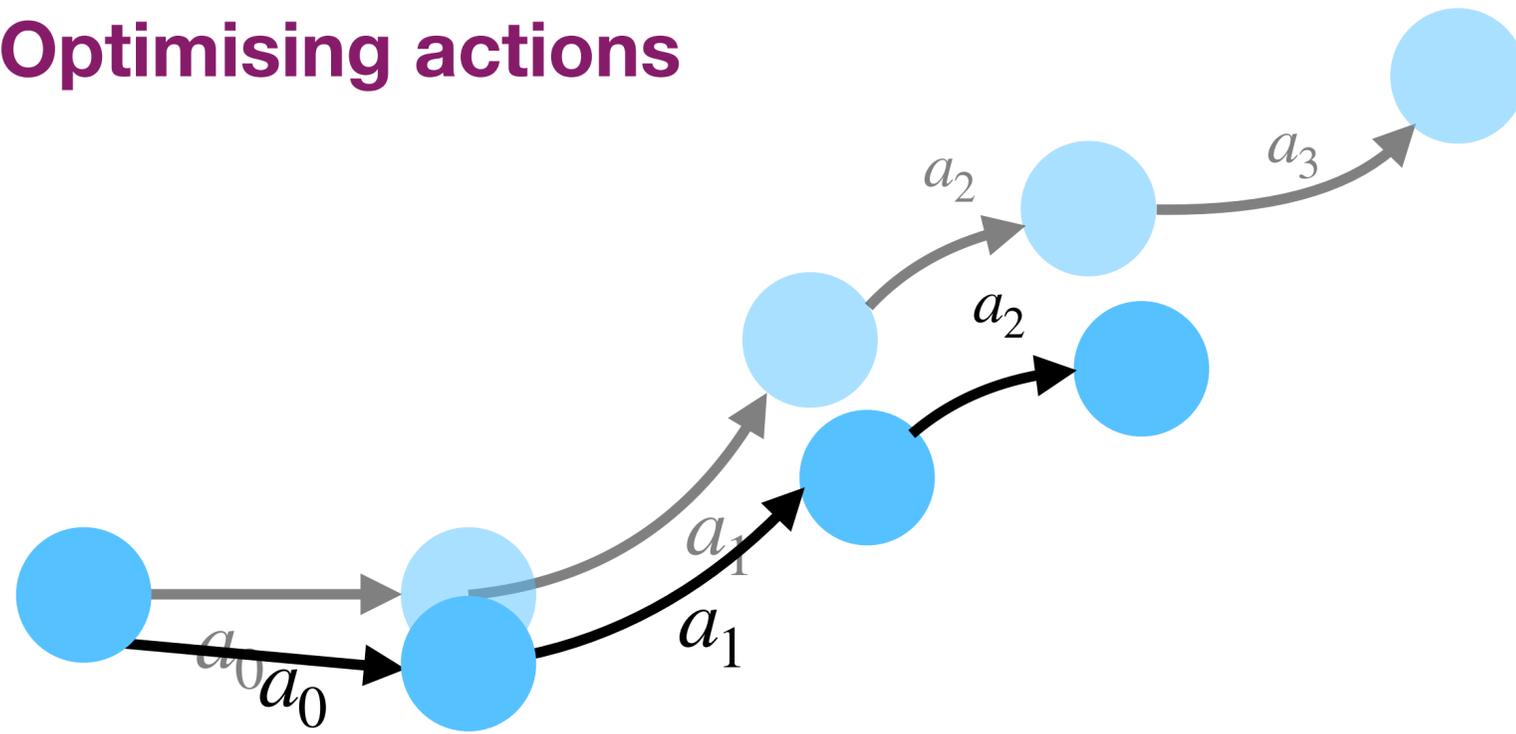
Recursively evaluate dynamics

# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



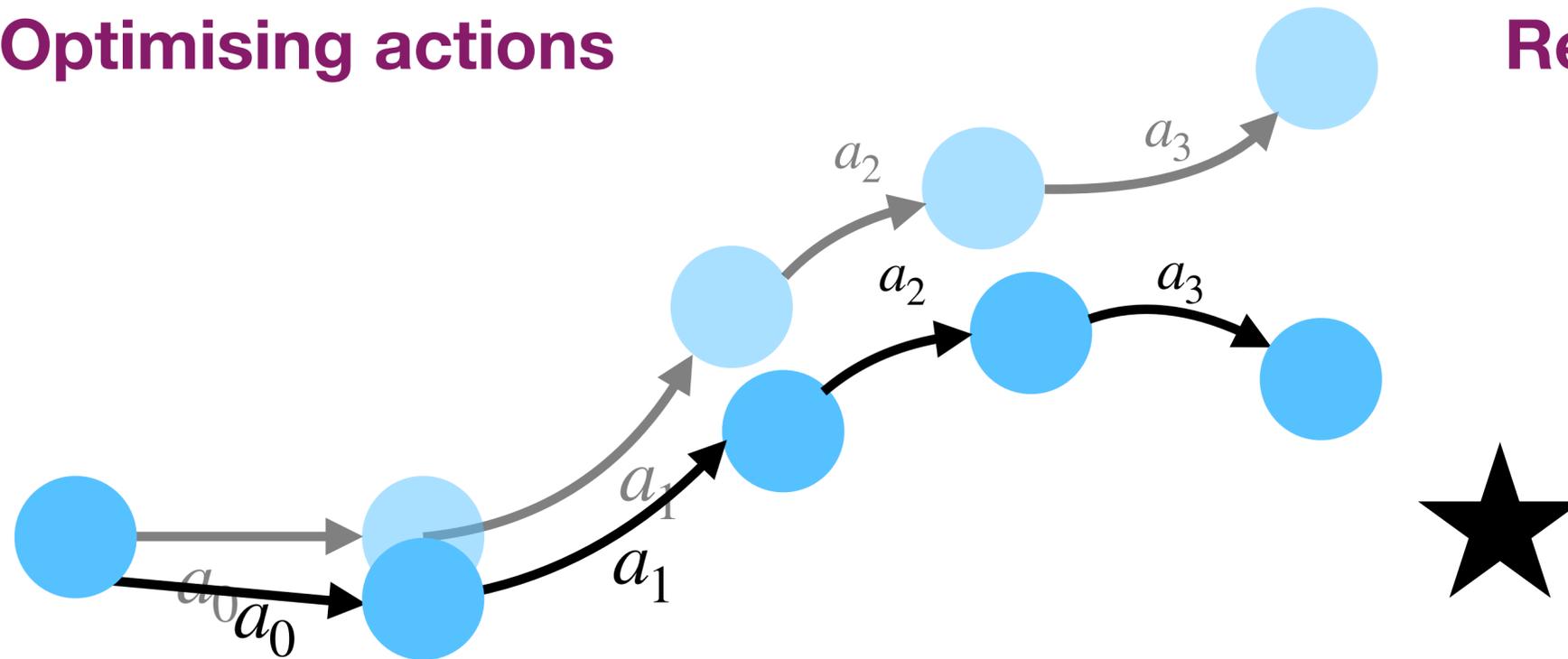
Recursively evaluate dynamics

# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



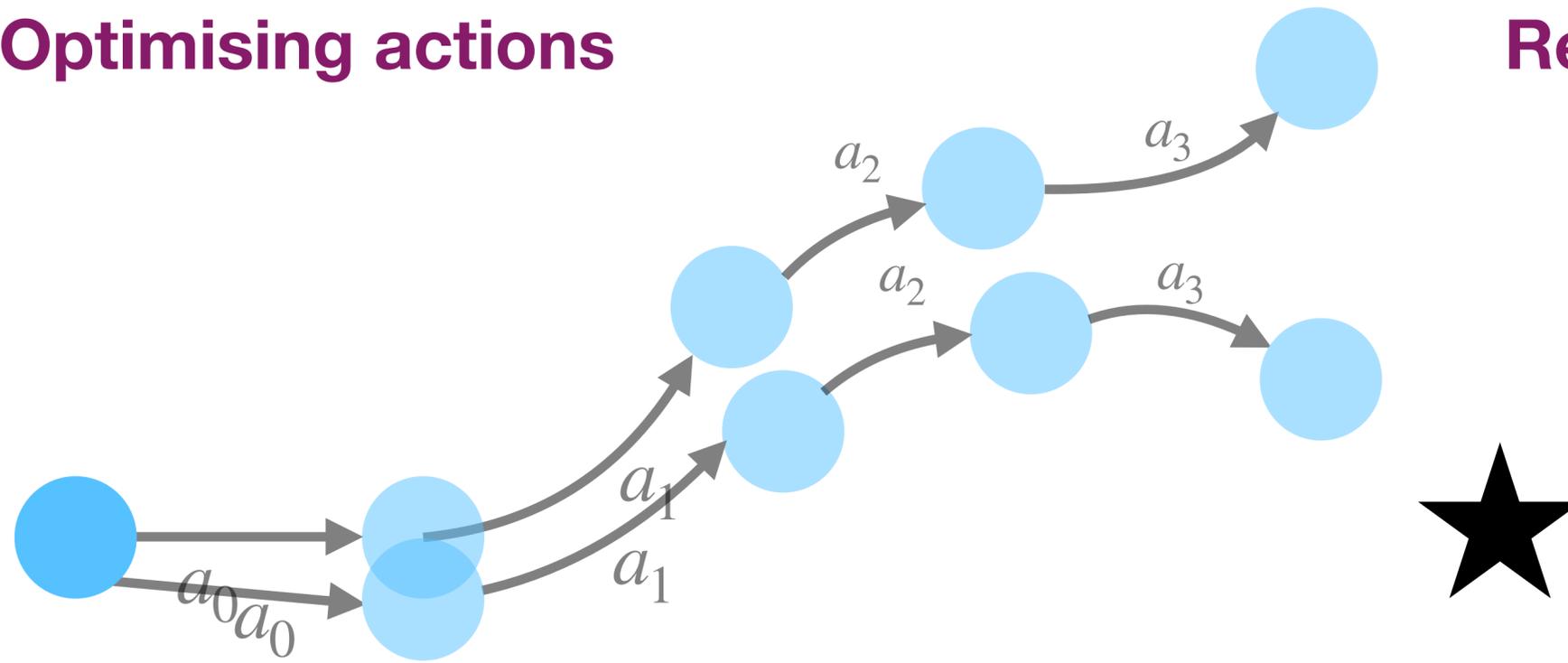
Recursively evaluate dynamics

# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



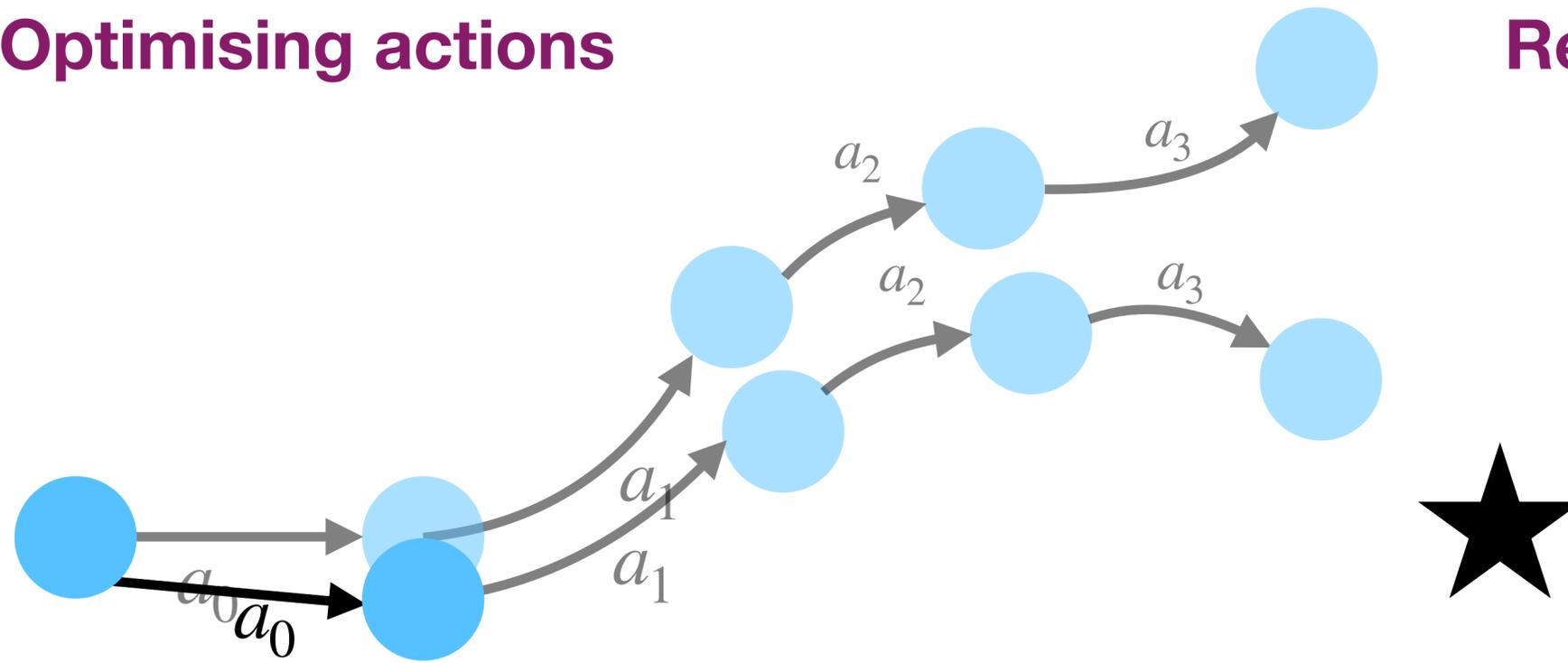
Recursively evaluate dynamics

# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



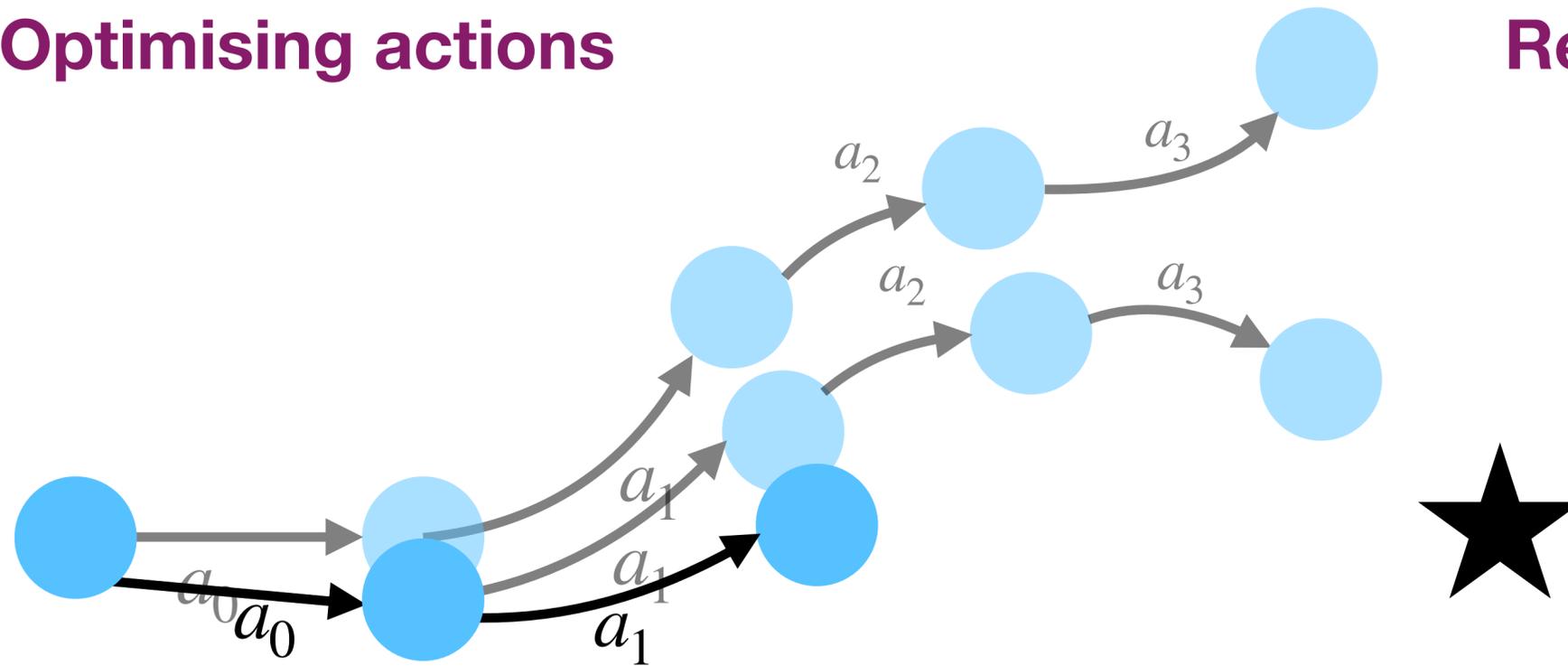
Recursively evaluate dynamics

# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



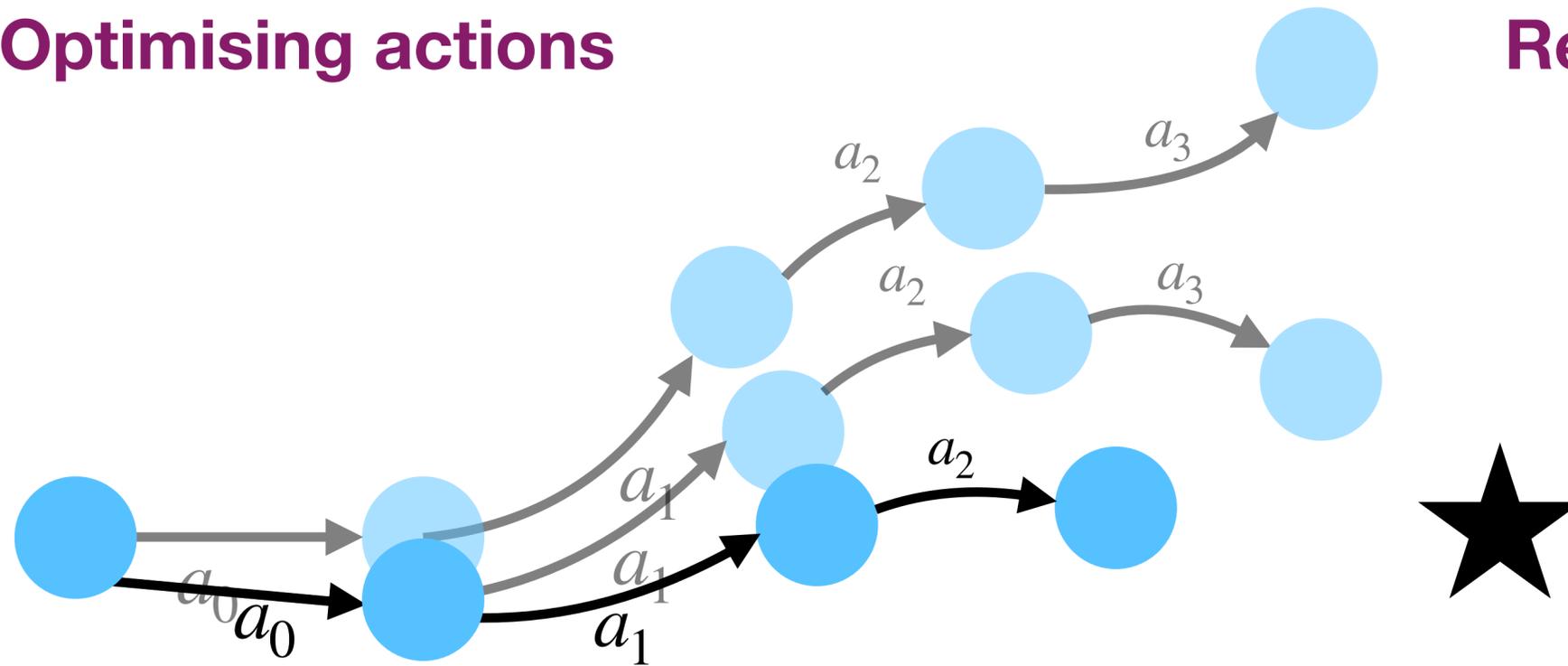
Recursively evaluate dynamics

# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



Recursively evaluate dynamics

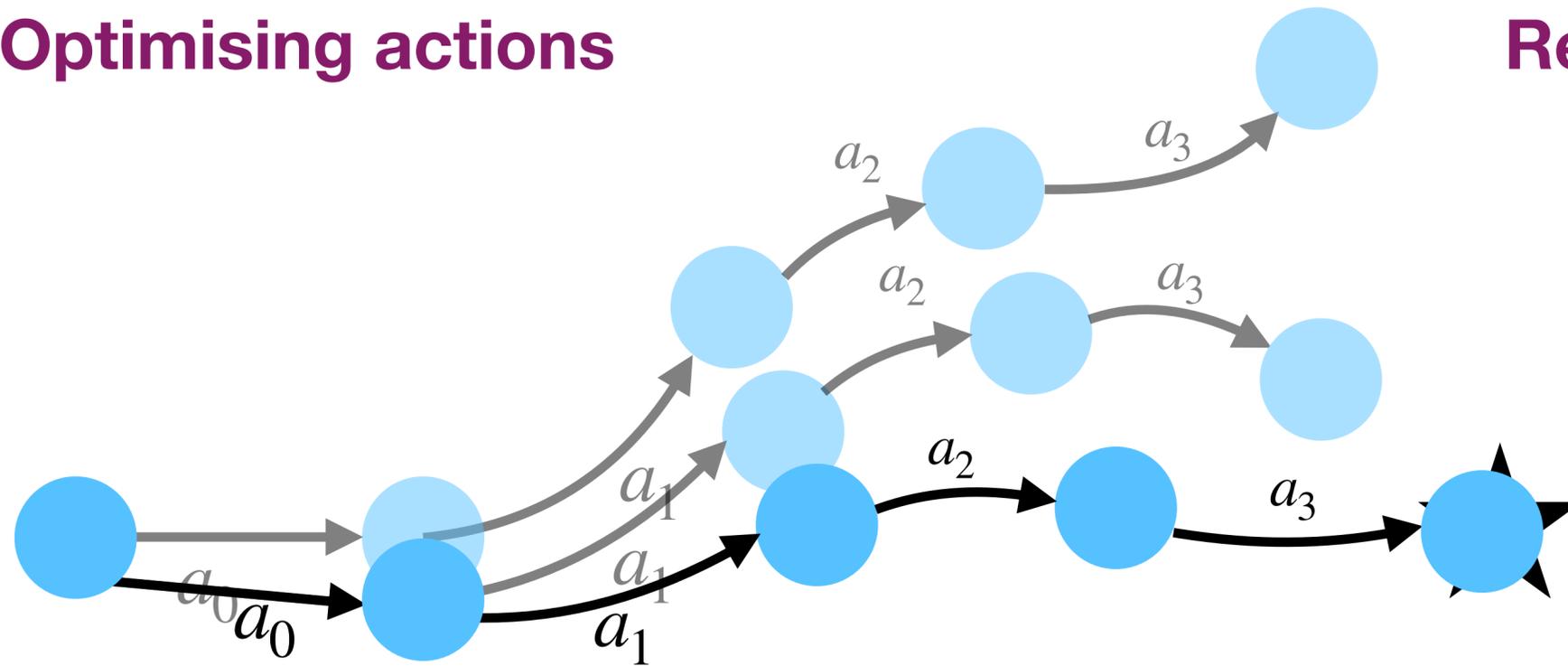
# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions

Recursively evaluate dynamics



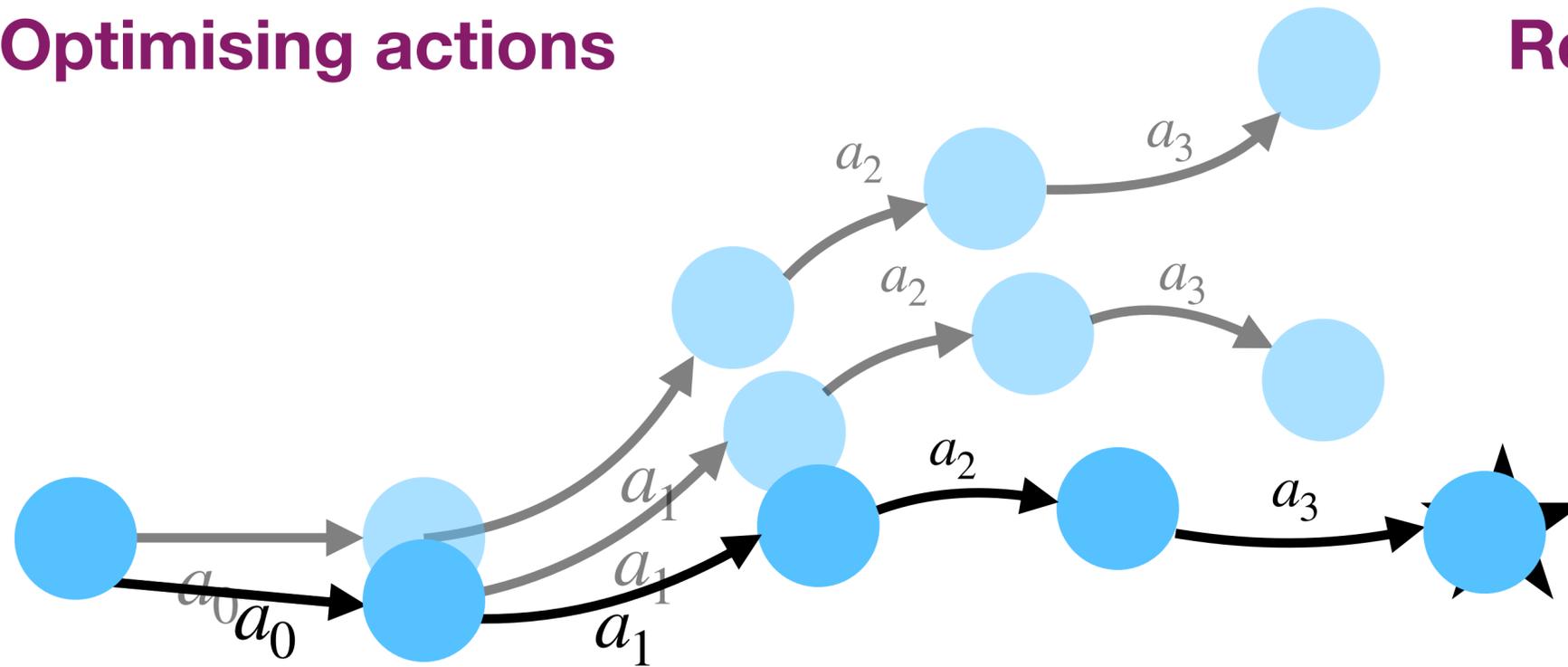
# Shooting Methods

## Illustration

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions

Recursively evaluate dynamics



# Shooting Methods

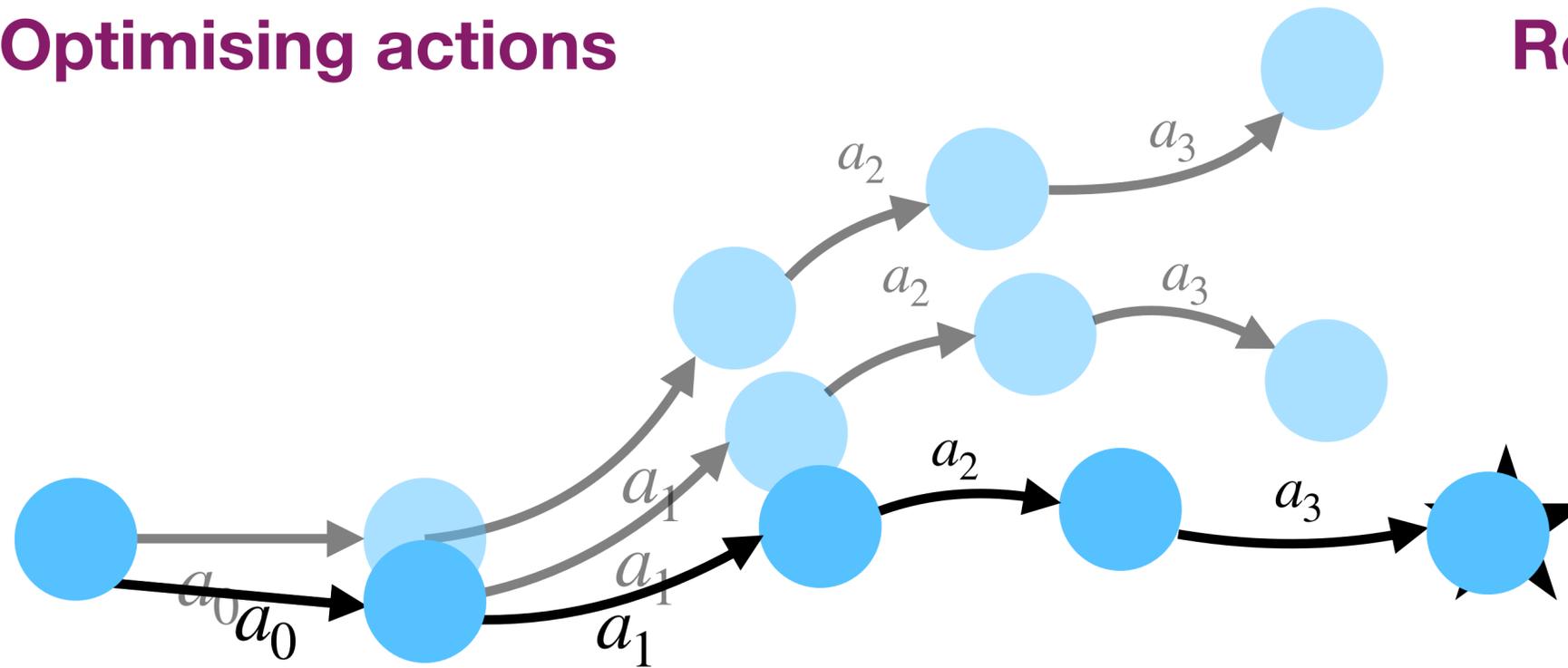
## Illustration

Gradient based approaches are fast  
But local minima

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions

Recursively evaluate dynamics



# Shooting Methods

## Illustration

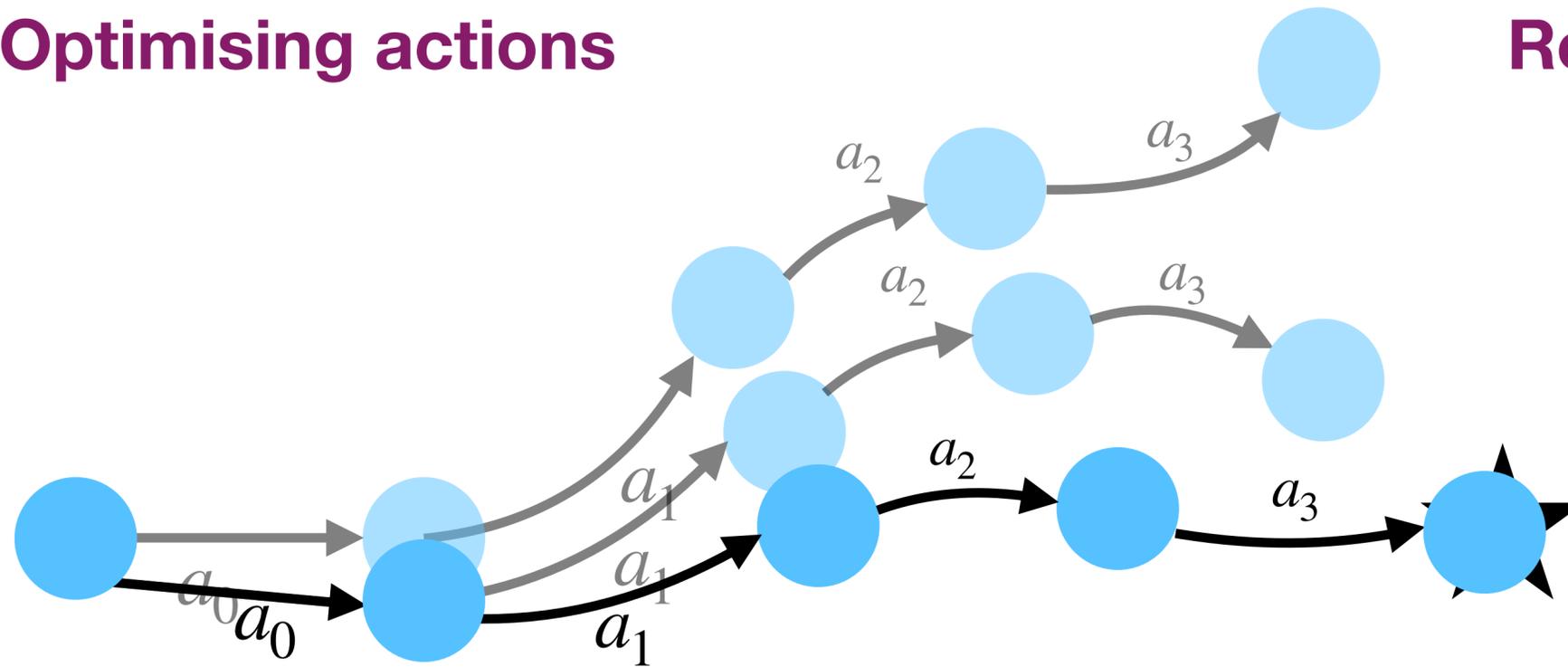
Gradient based approaches are fast

But local minima

And vanishing/exploding gradients

$$J(a_{0:H}) = \gamma^0 r(s_0, a_0) + \gamma^1 r(f(s_0, a_0), a_1) + \dots + \gamma^H r(f(f(\dots), a_{H-1}), a_H)$$

Optimising actions



Recursively evaluate dynamics

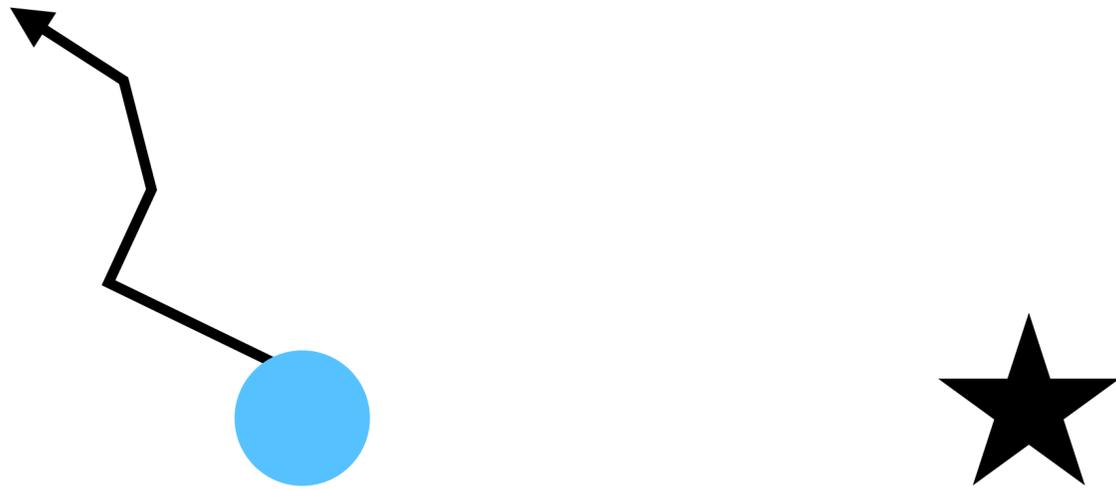
# Shooting Methods

## Random shooting



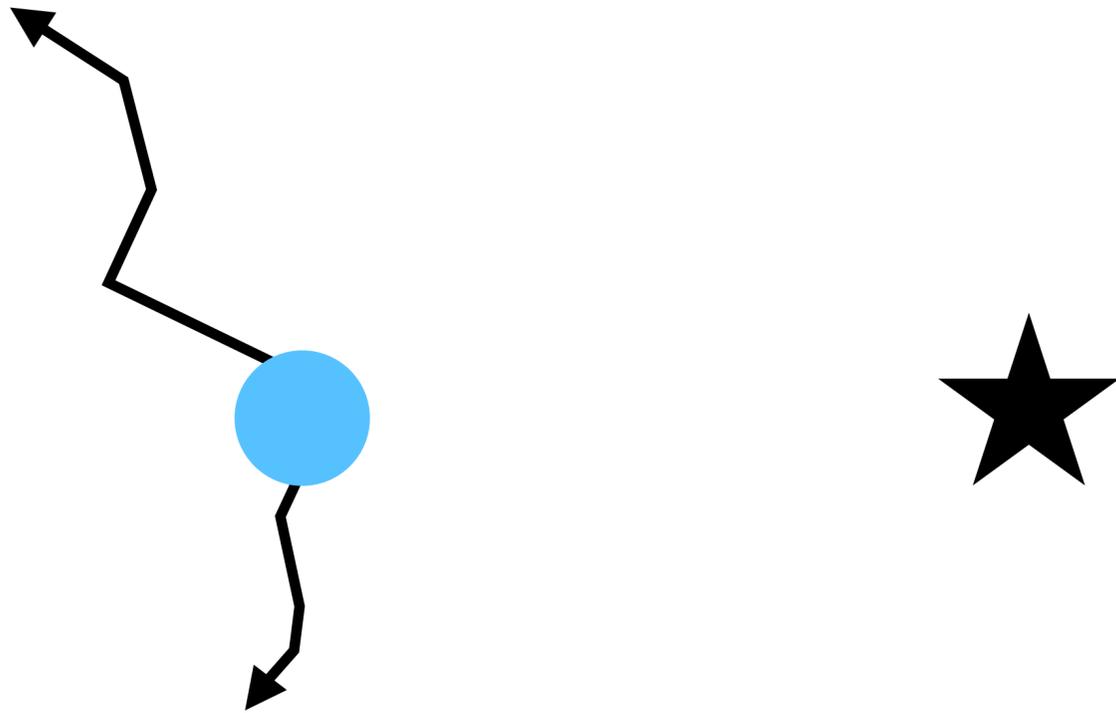
# Shooting Methods

## Random shooting



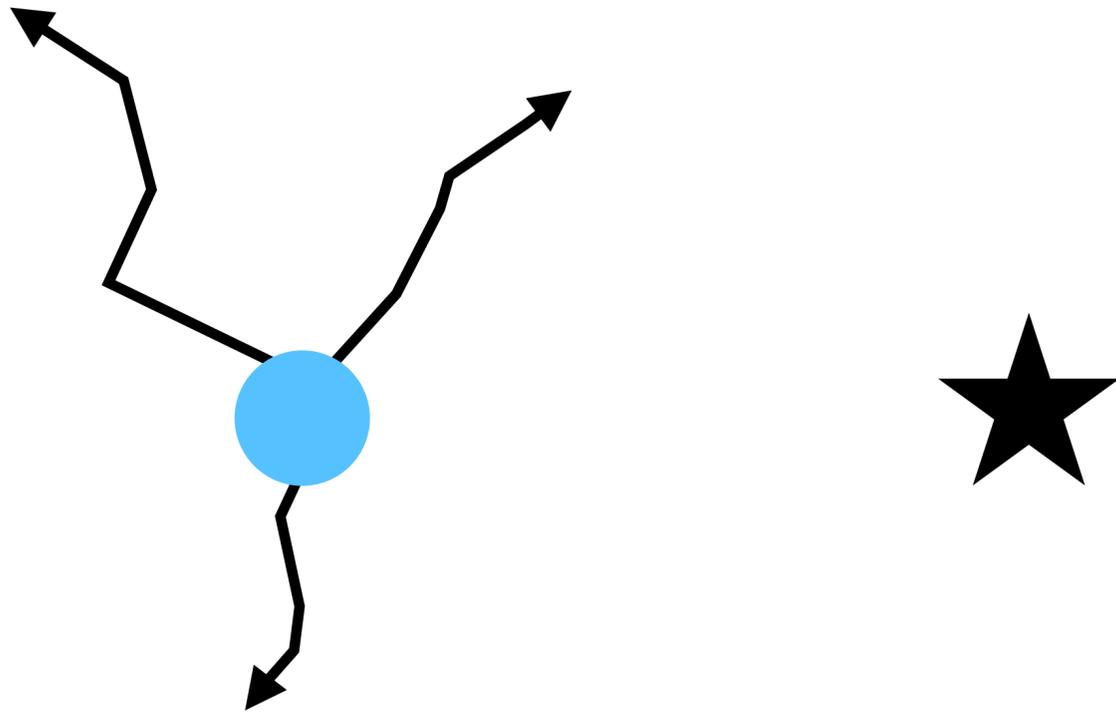
# Shooting Methods

## Random shooting



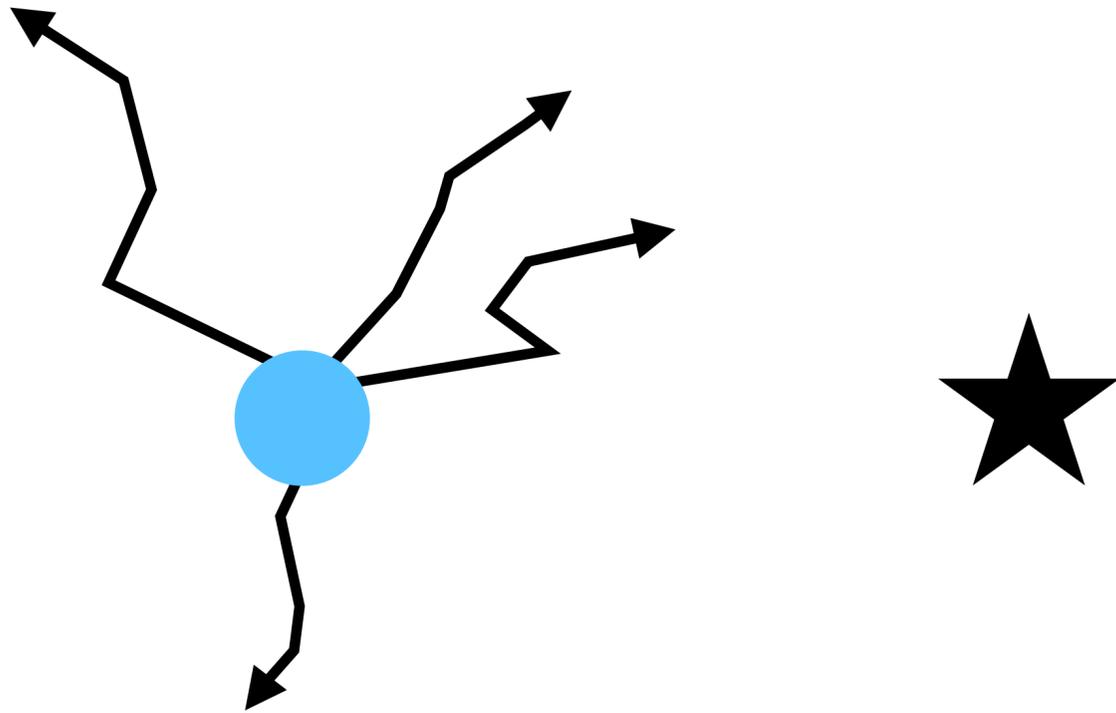
# Shooting Methods

## Random shooting



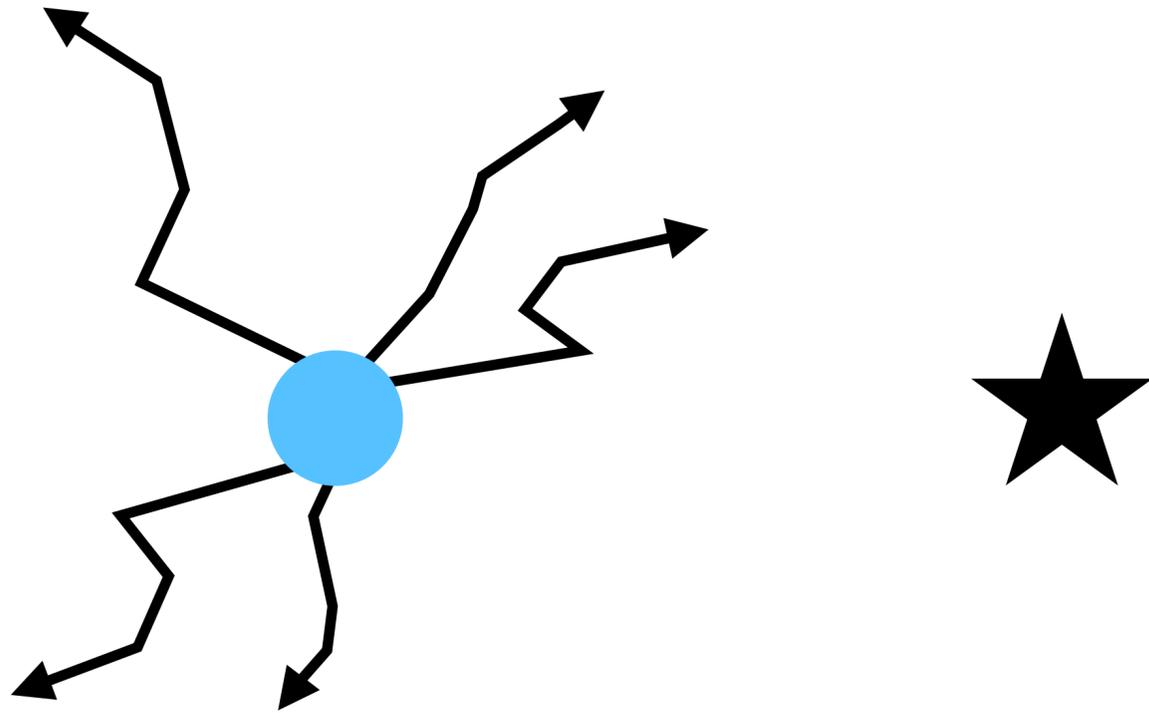
# Shooting Methods

## Random shooting



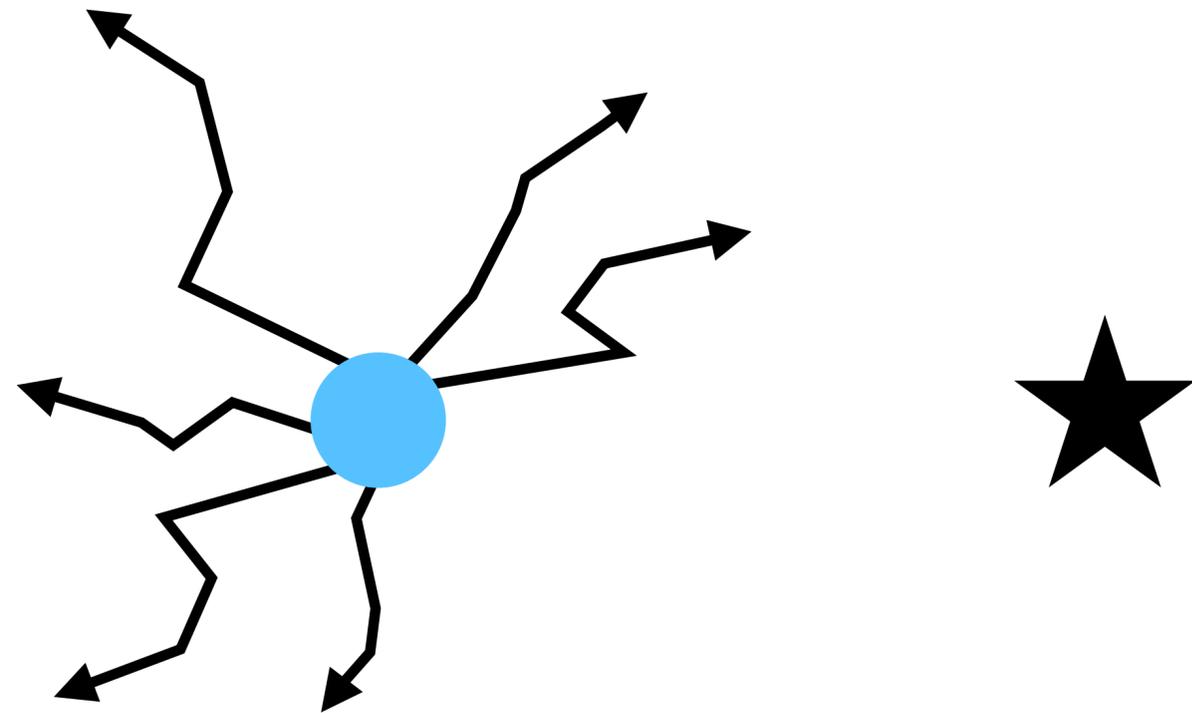
# Shooting Methods

## Random shooting



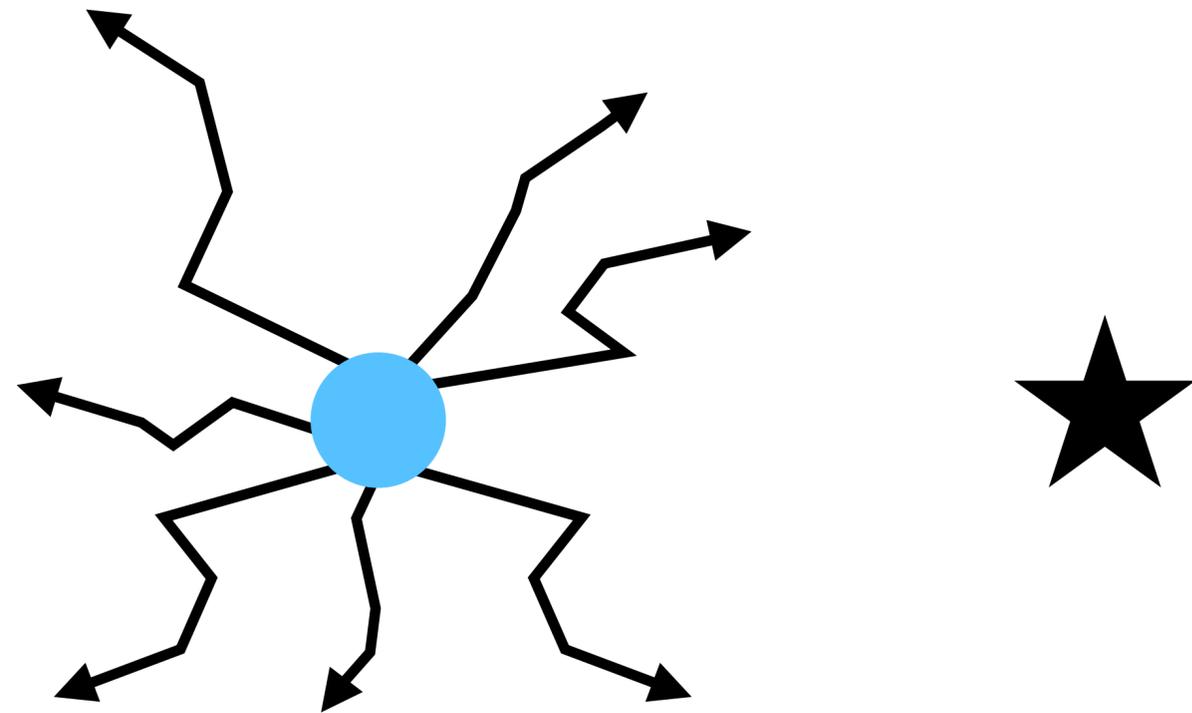
# Shooting Methods

## Random shooting



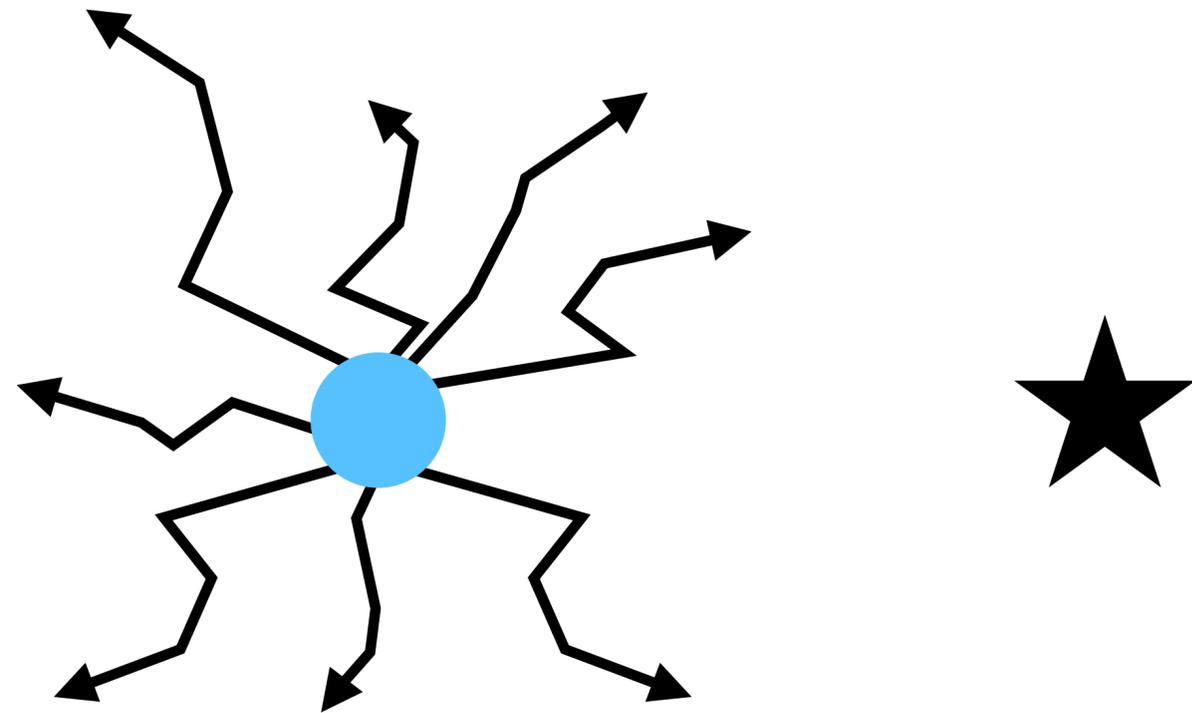
# Shooting Methods

## Random shooting



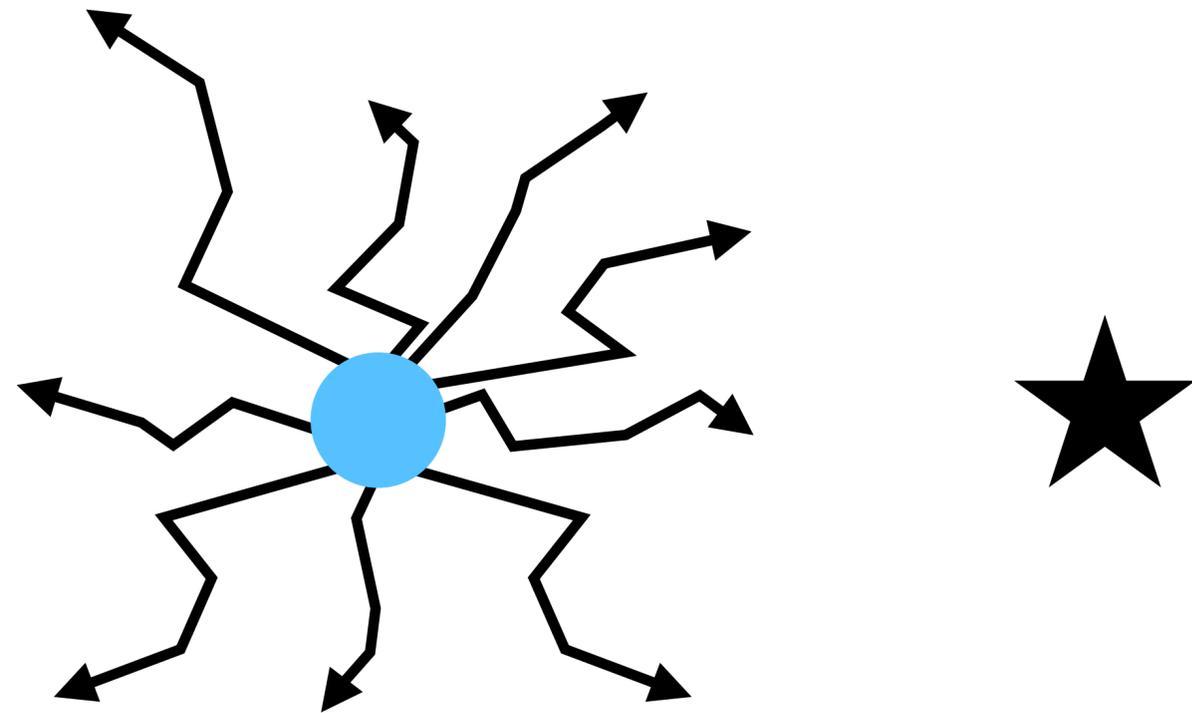
# Shooting Methods

## Random shooting



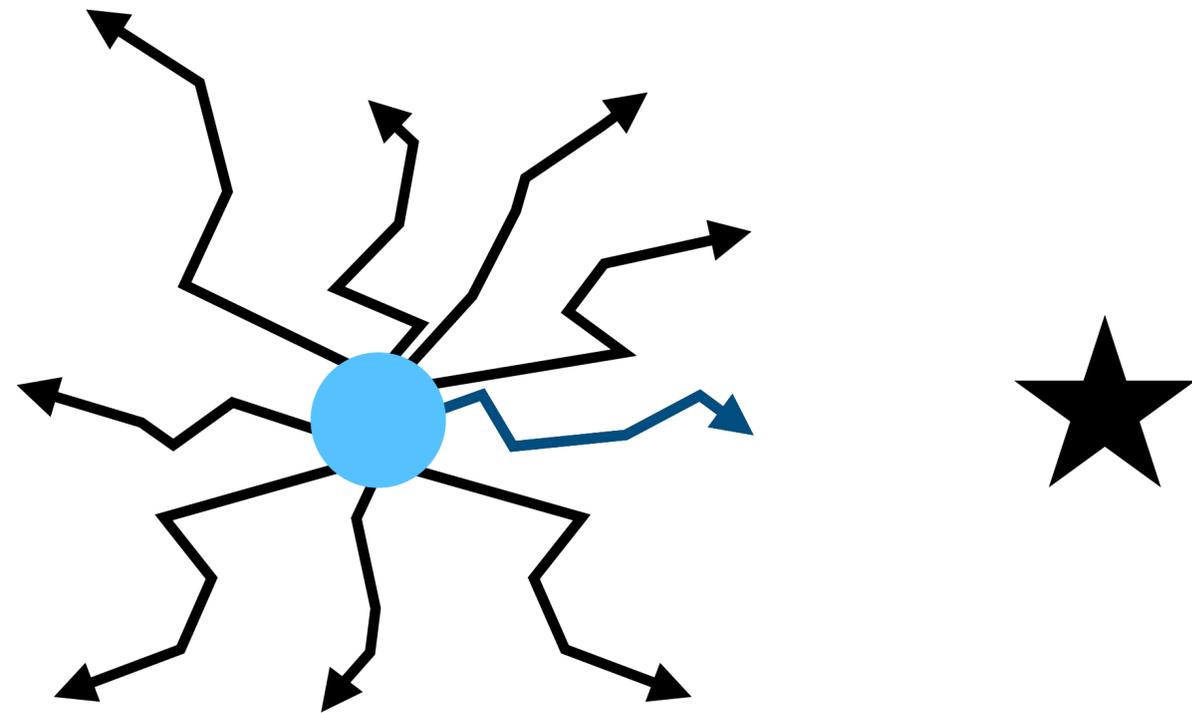
# Shooting Methods

## Random shooting



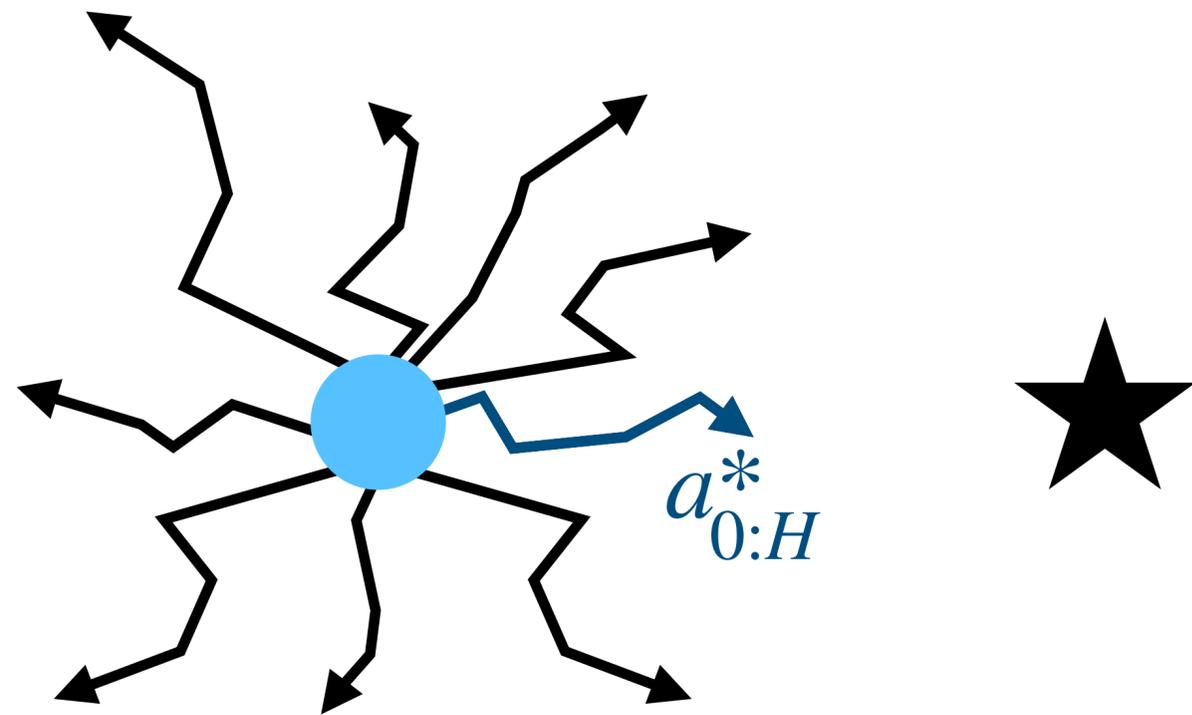
# Shooting Methods

## Random shooting



# Shooting Methods

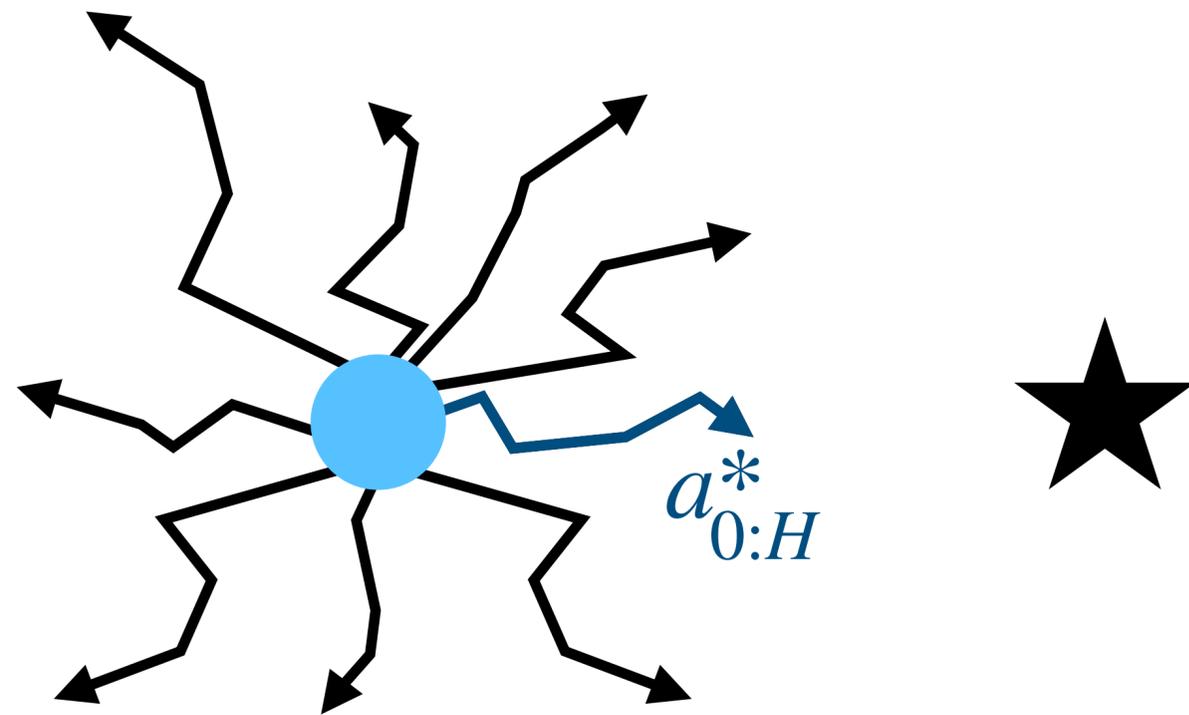
## Random shooting



# Shooting Methods

## Random shooting

Simple



FCAI

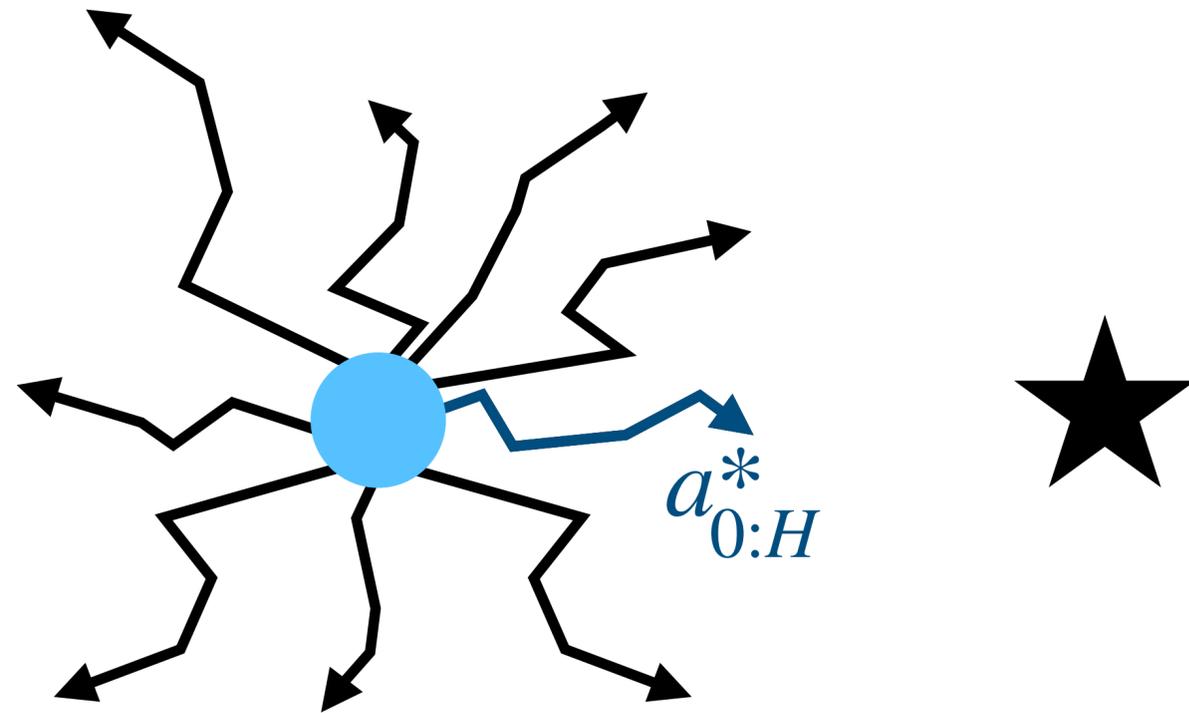
fcai.fi

# Shooting Methods

Random shooting

Simple

Parallelisable



FCAI

fcai.fi

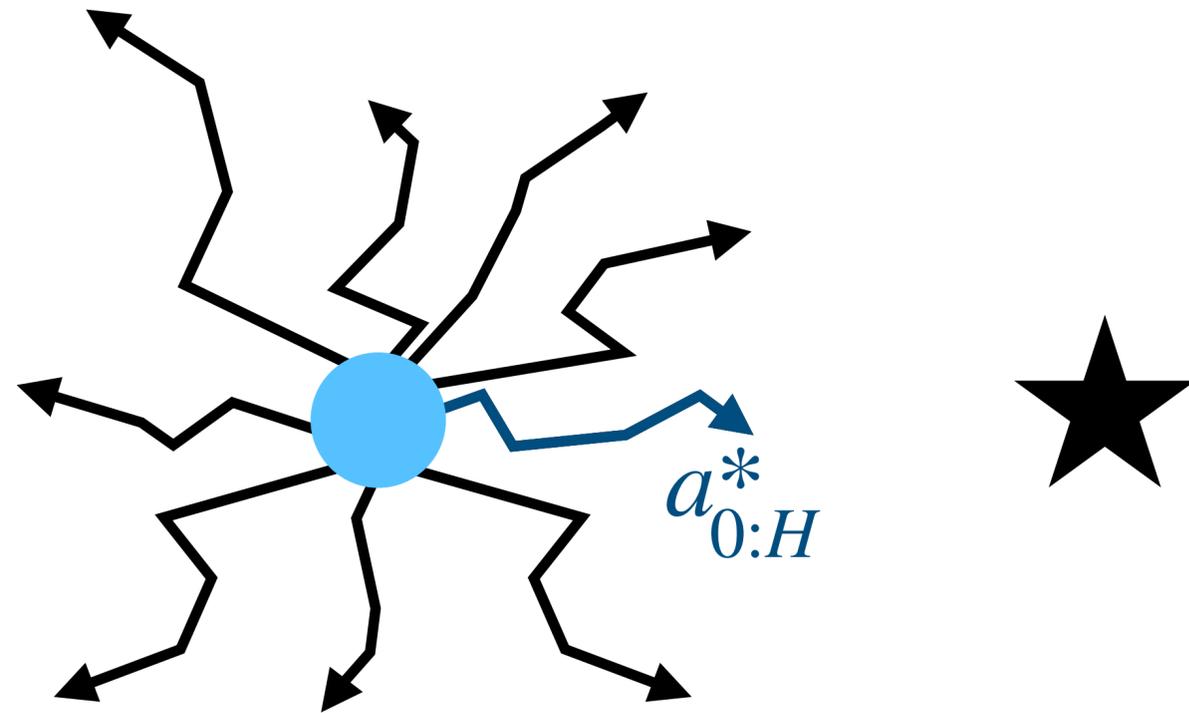
# Shooting Methods

Random shooting

Simple

Parallelisable

Sample inefficient



**FCAI**

**fcai.fi**

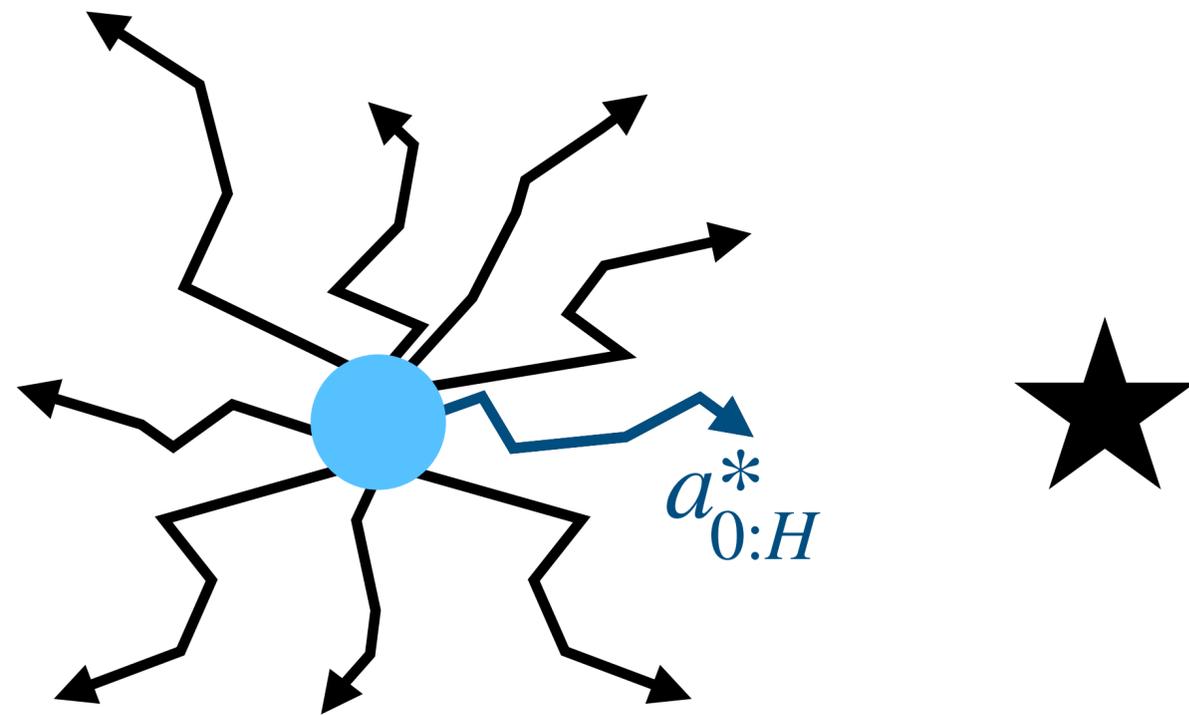
# Shooting Methods

Random shooting

Simple

Parallelisable

Sample inefficient



**FCAI**

**fcai.fi**

# Shooting Methods

## Cross-Entropy Method

Iteration 1



# Shooting Methods

## Cross-Entropy Method

Iteration 1

**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$



# Shooting Methods

## Cross-Entropy Method

Iteration 1

**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration



# Shooting Methods

## Cross-Entropy Method

Iteration 1

**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

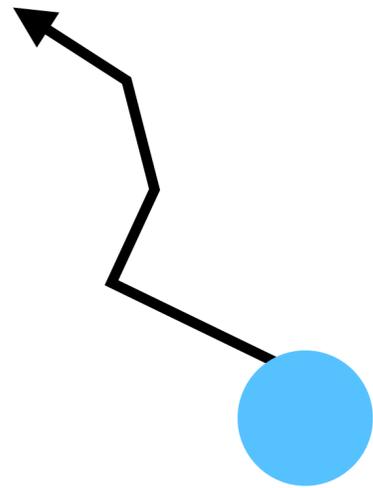
**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution



# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

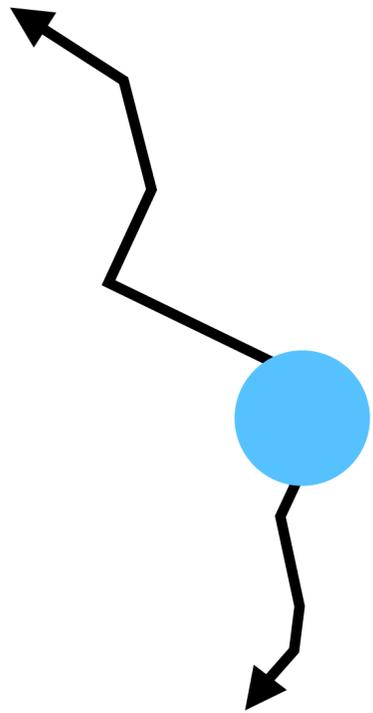
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

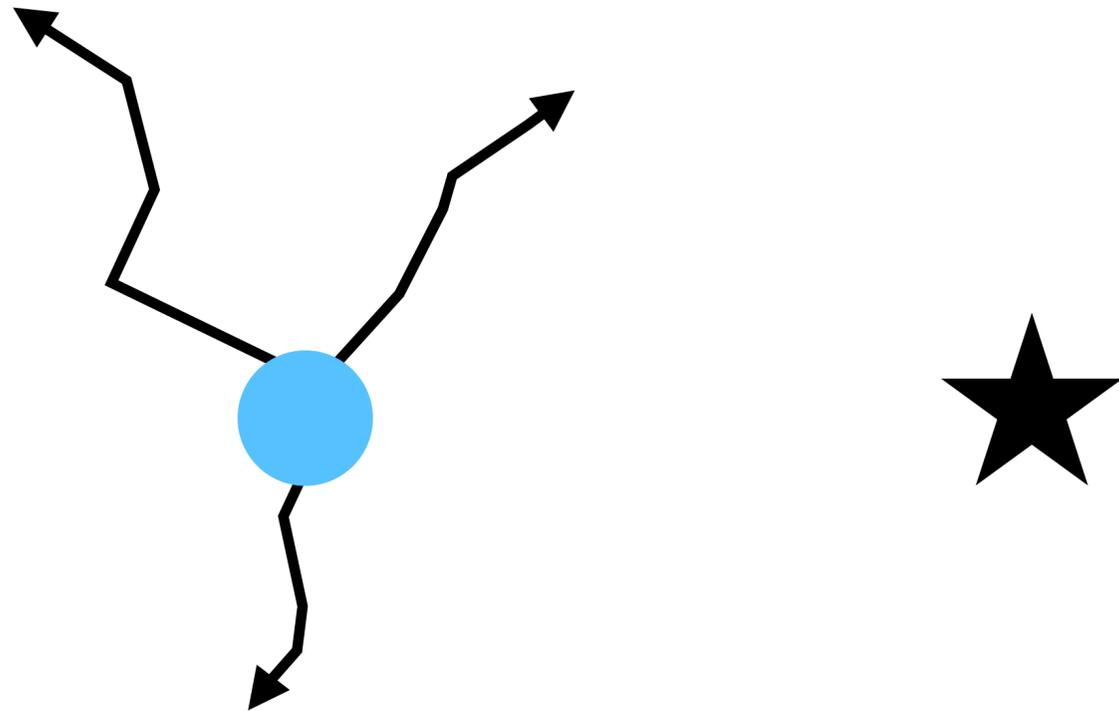
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

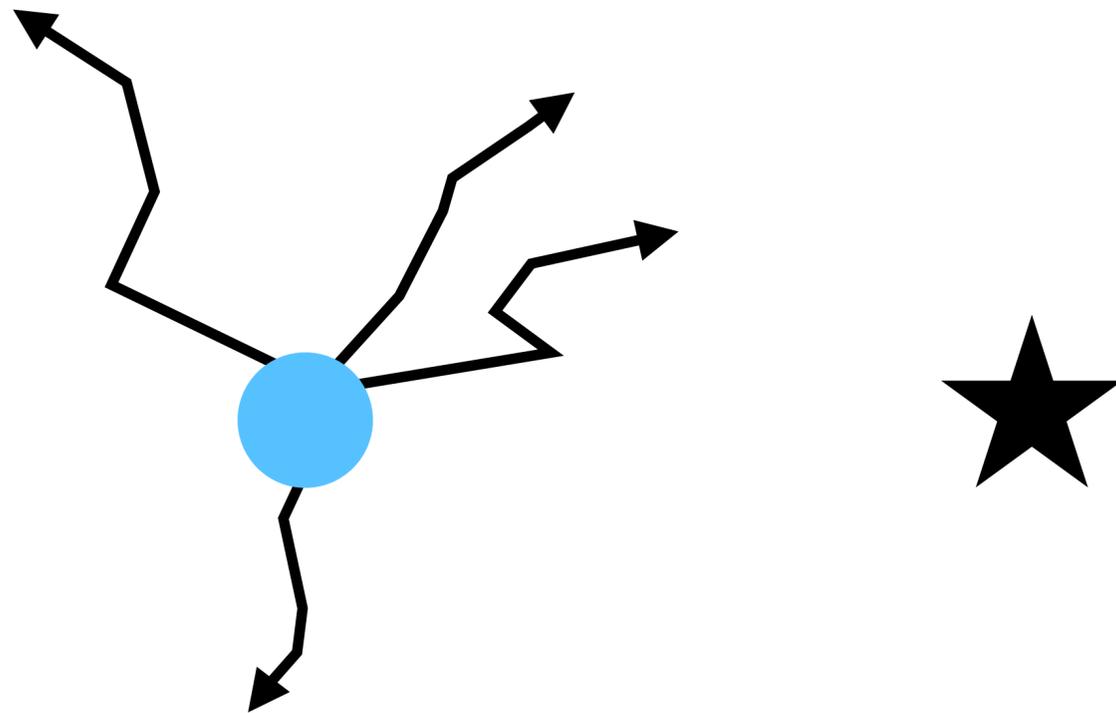
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

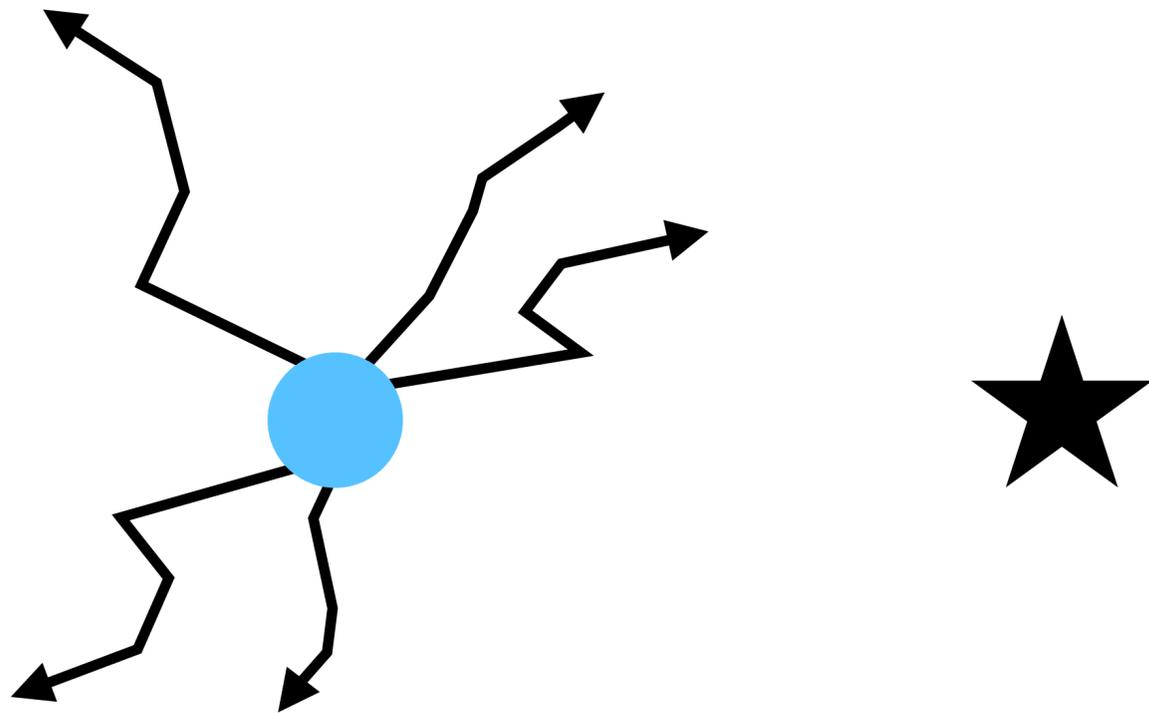
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

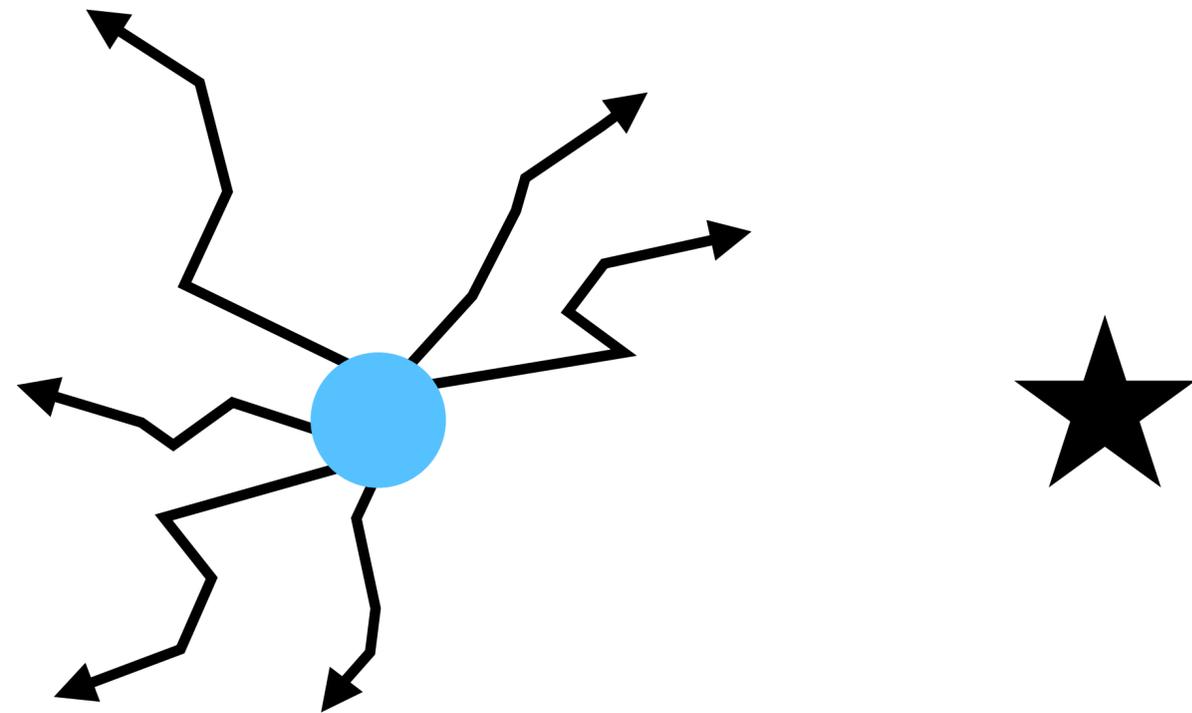
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

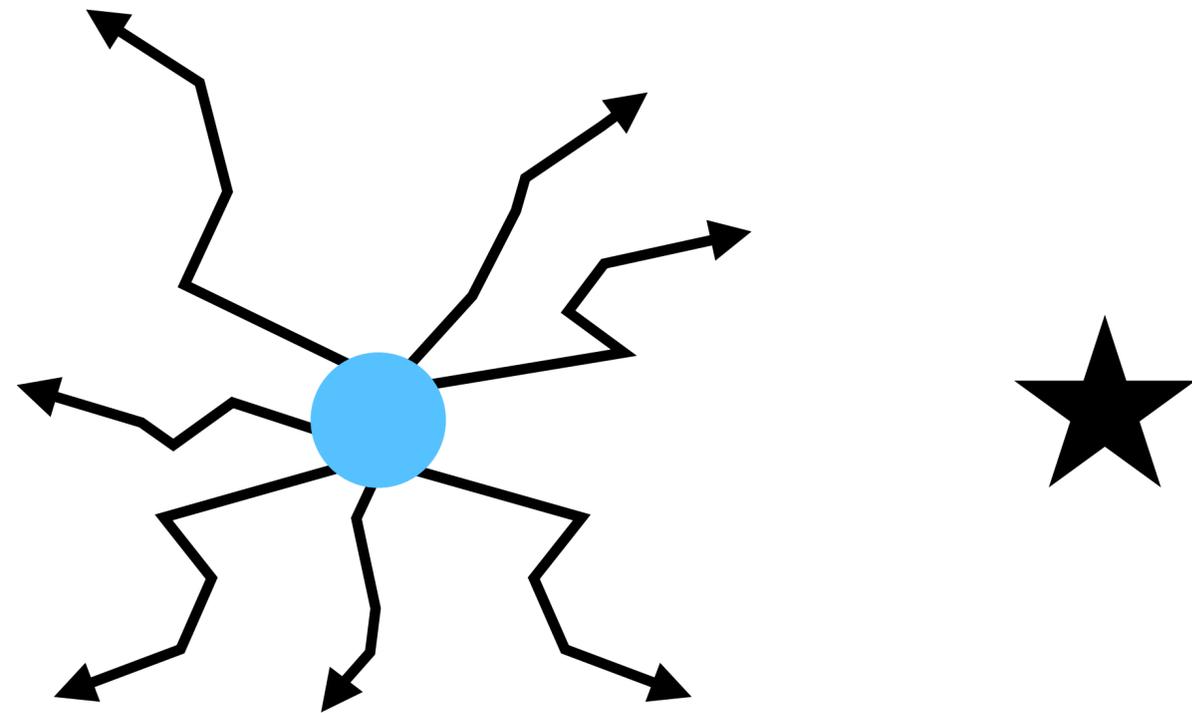
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

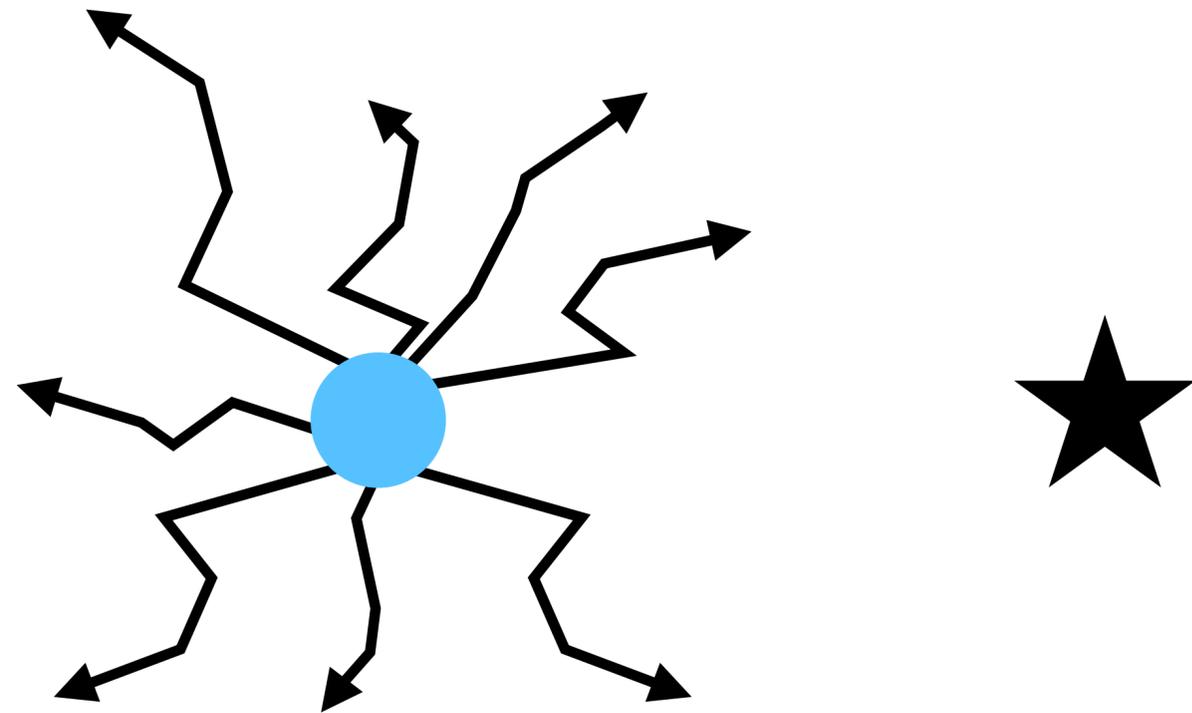
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

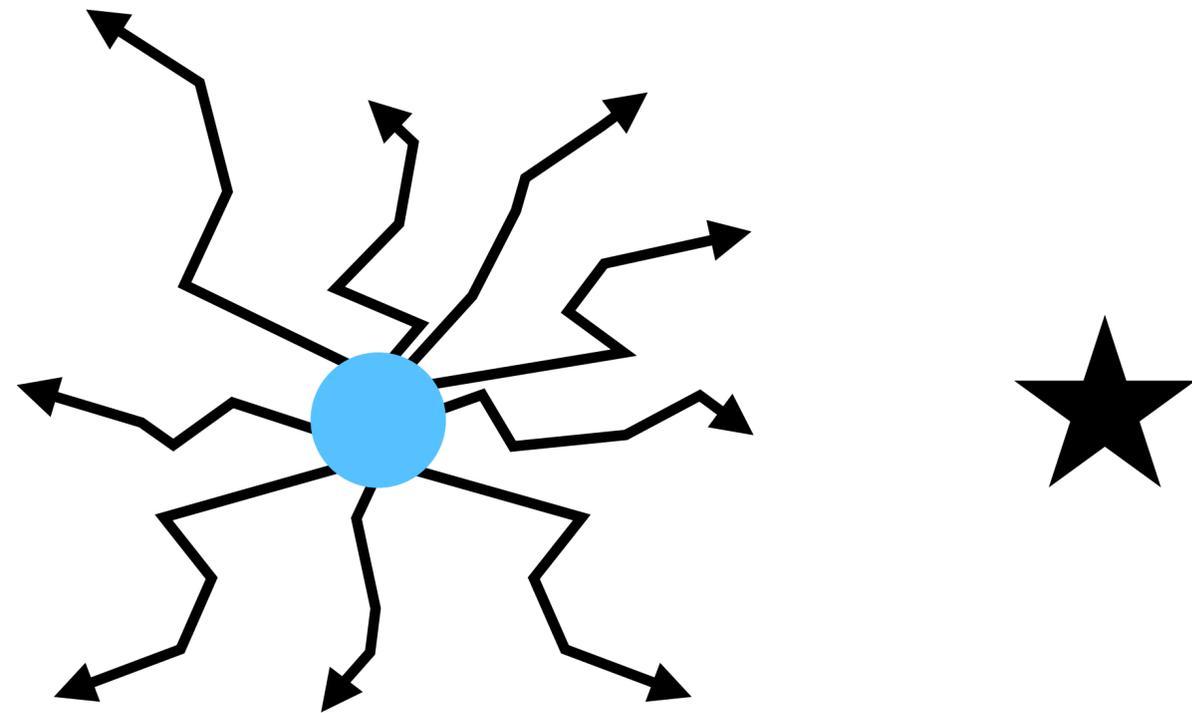
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

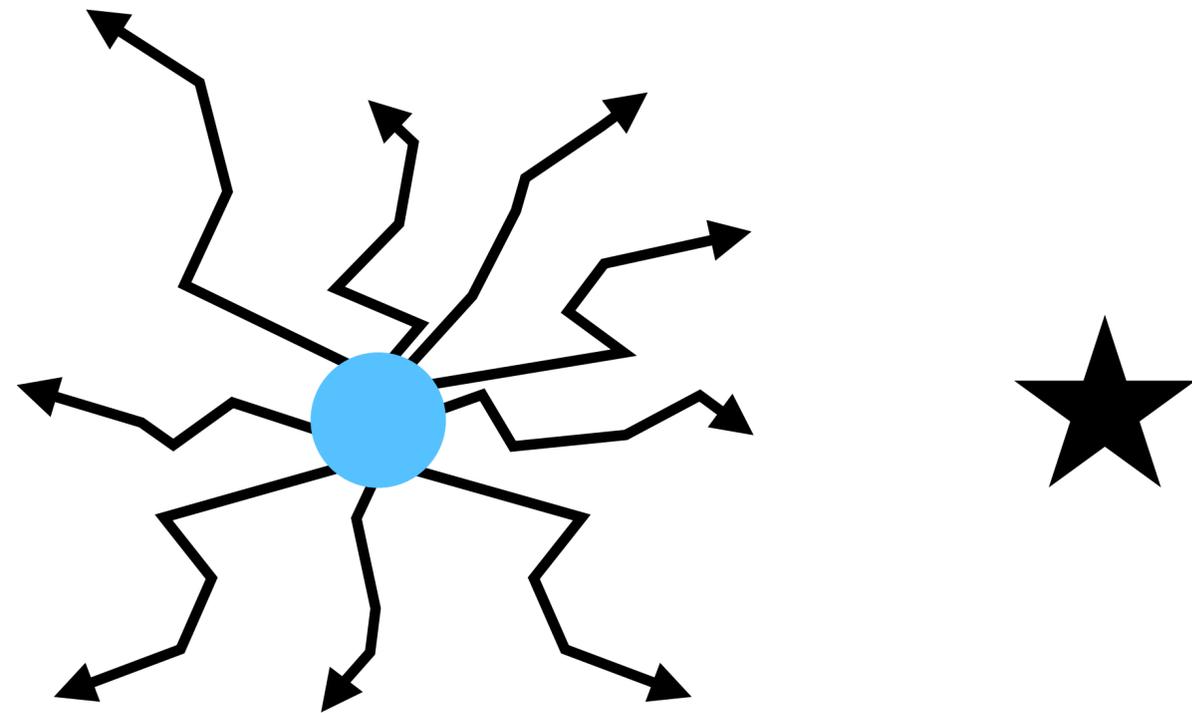
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

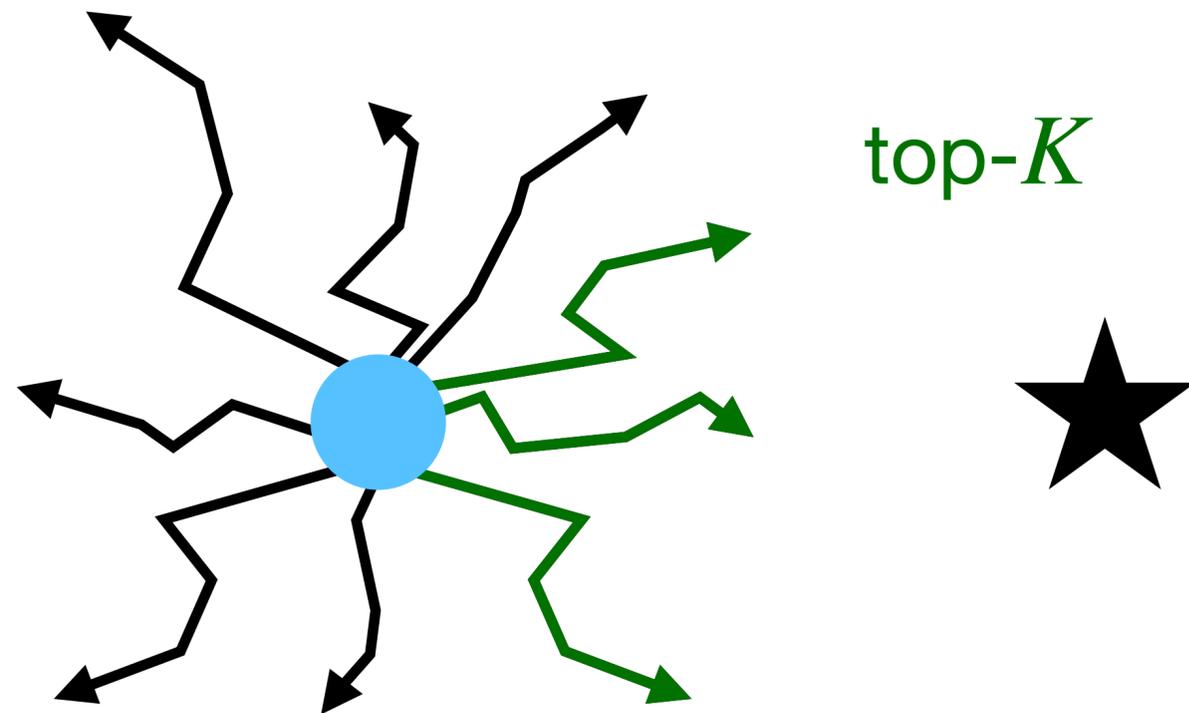
**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

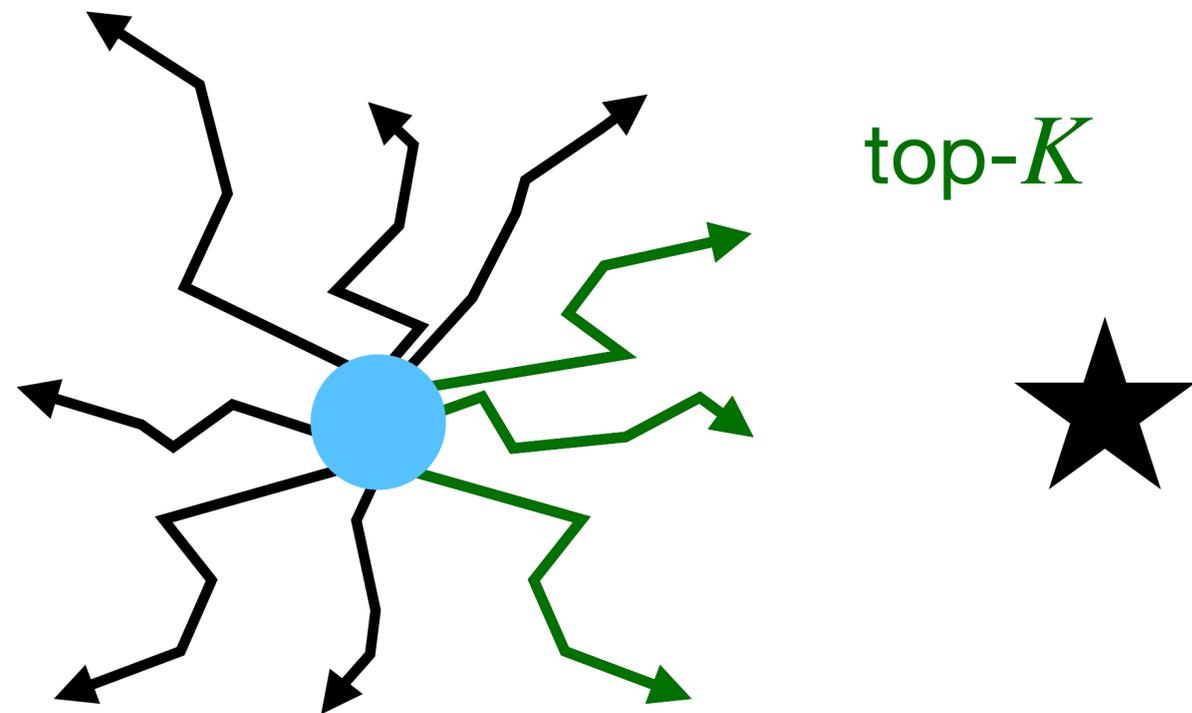
**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

# Shooting Methods

## Cross-Entropy Method

Iteration 1



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

### Iteration 2

**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

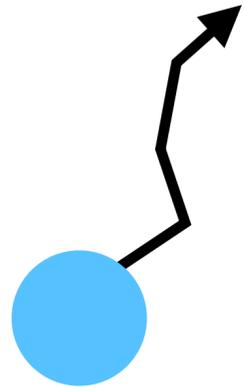
**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples



# Shooting Methods

## Cross-Entropy Method

Iteration 2



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

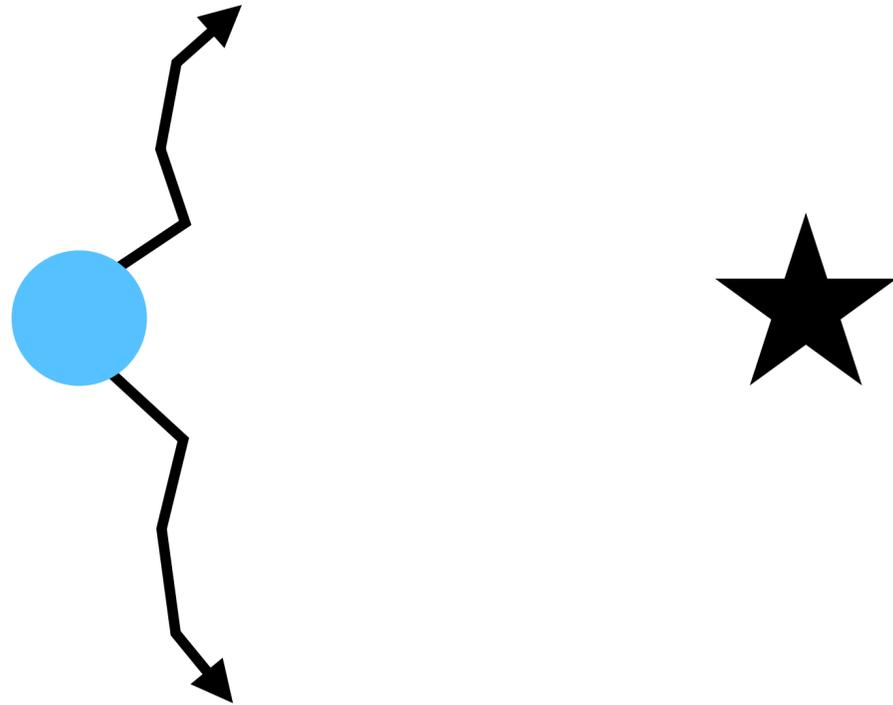
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 2



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

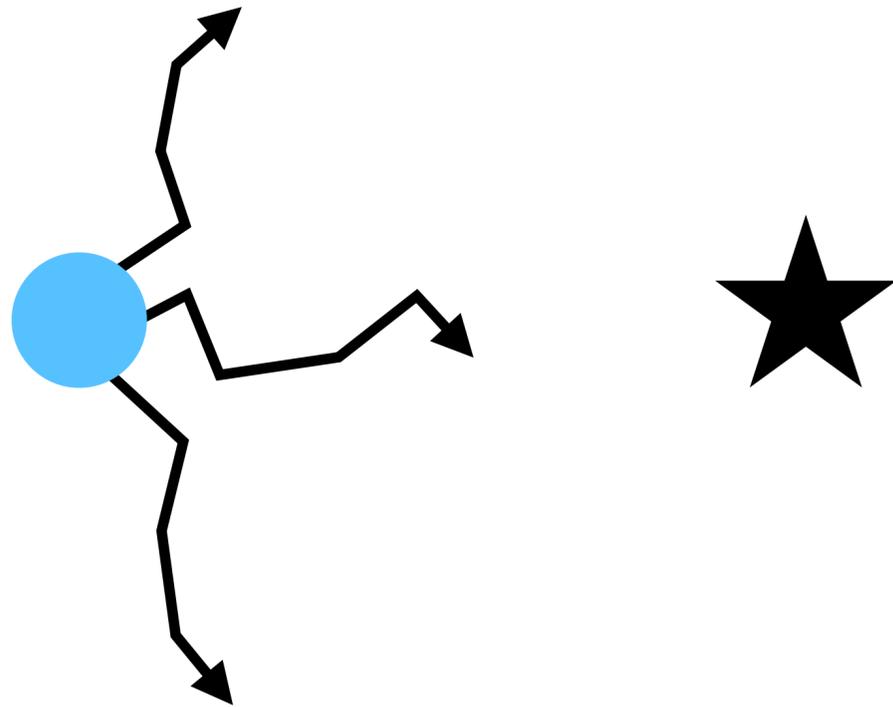
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 2



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

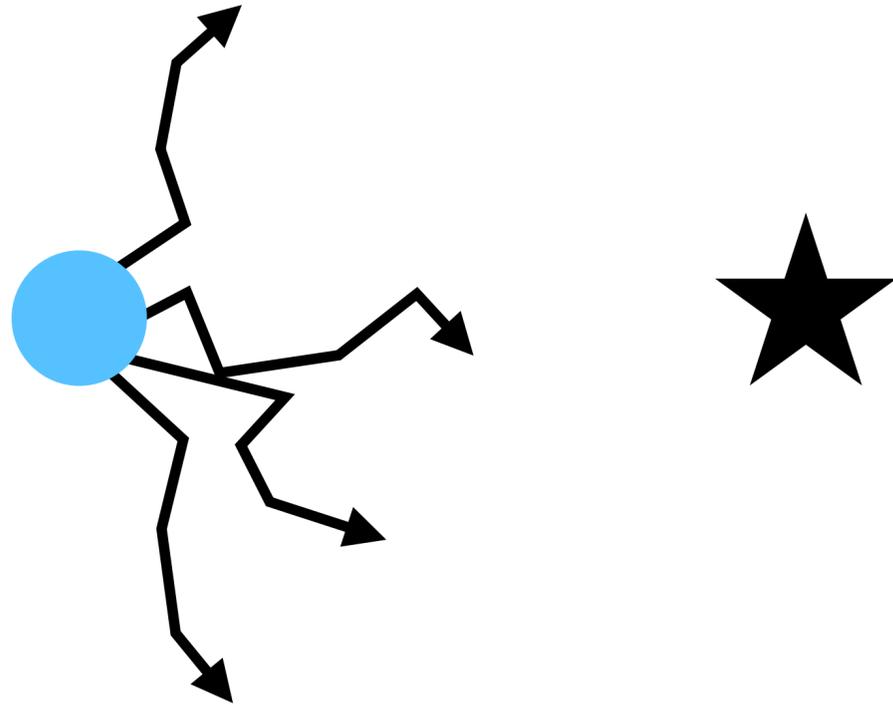
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 2



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

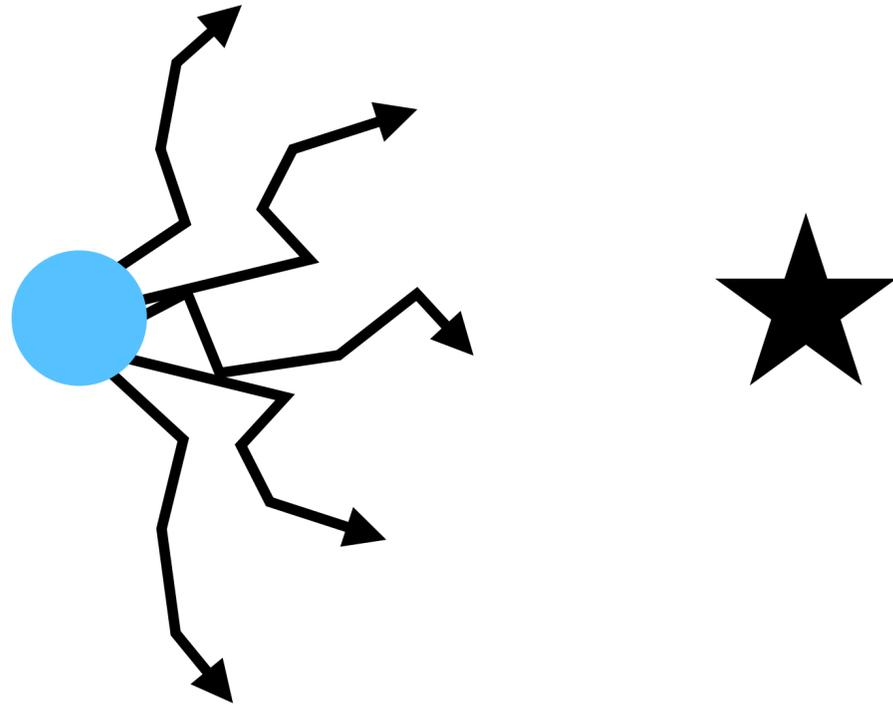
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 2



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

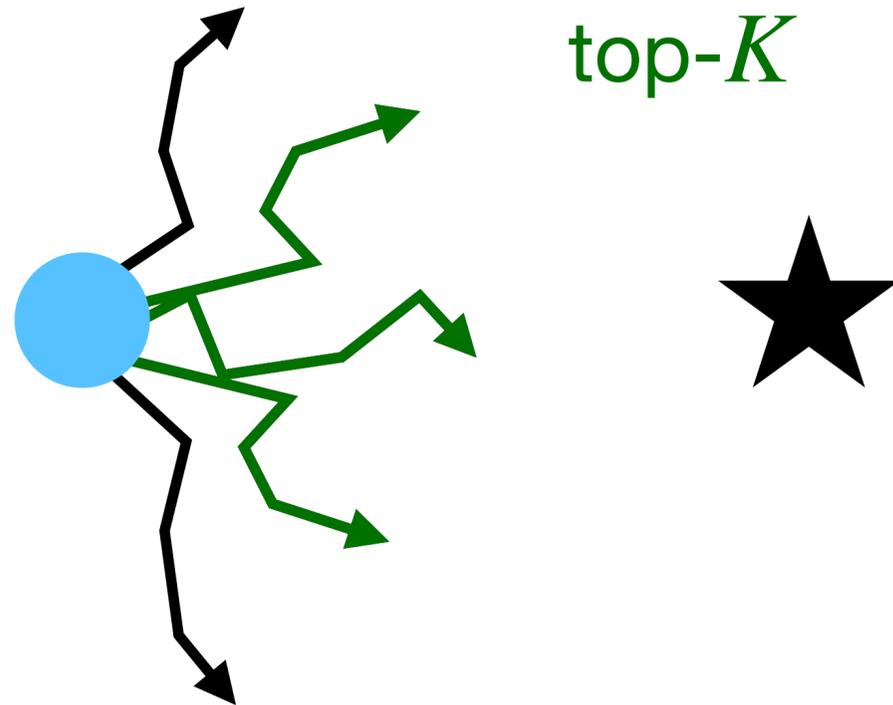
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 2



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

### Iteration 3

**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

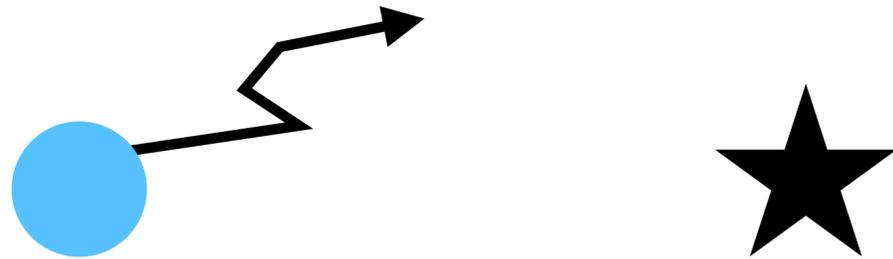
**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples



# Shooting Methods

## Cross-Entropy Method

Iteration 3



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

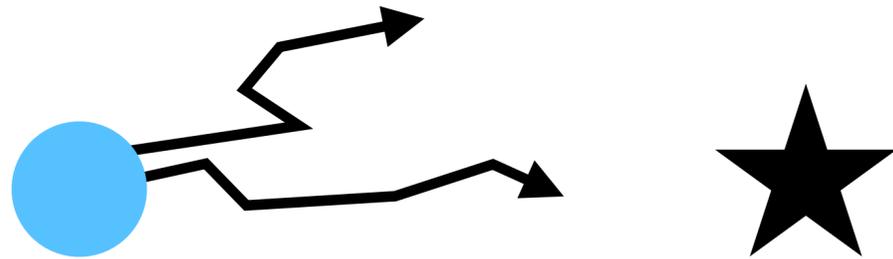
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 3



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

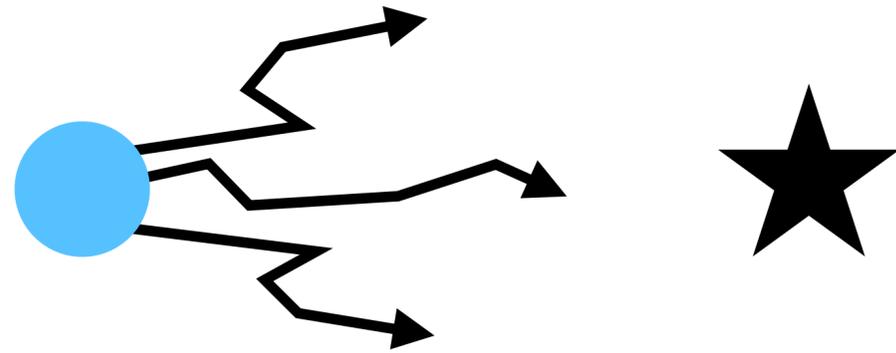
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 3



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

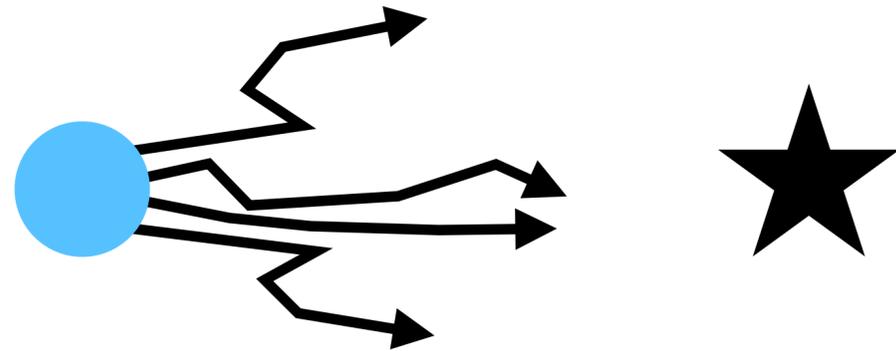
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 3



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

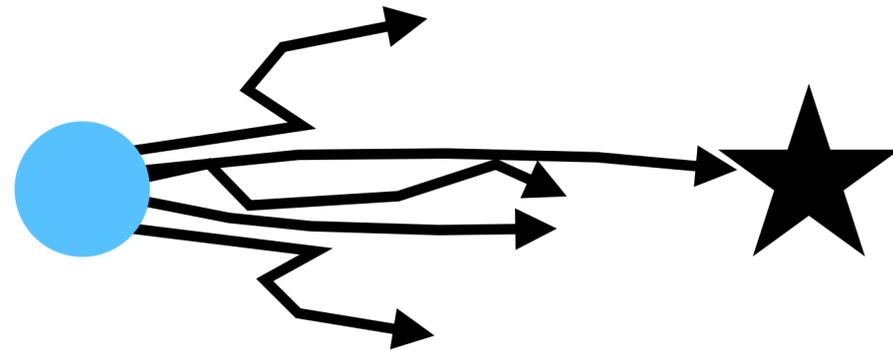
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 3



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

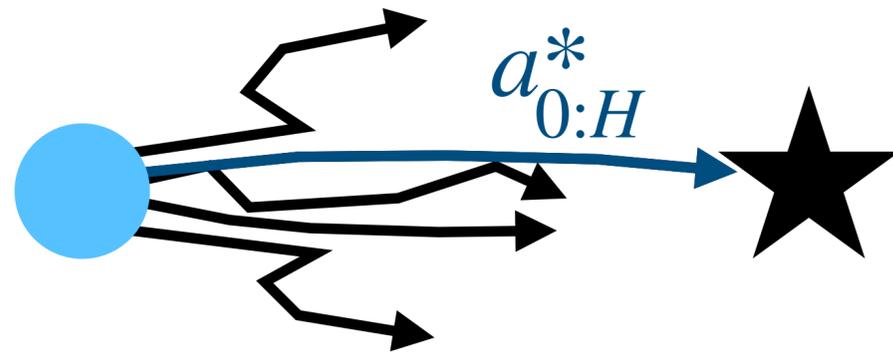
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 3



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

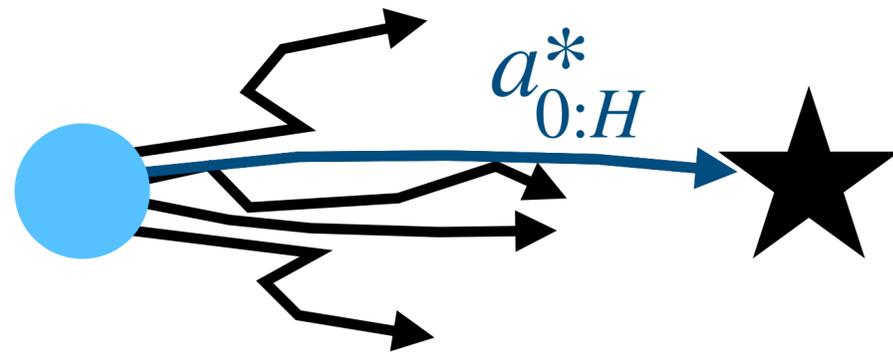
**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

Iteration 3



**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples

# Shooting Methods

## Cross-Entropy Method

More sample efficient

Faster convergence

Iteration 3

**Initialise** action sequence sampling distribution  $\{a_t \sim \mathcal{N}(\mu_t, \sigma_t^2)\}_{t=0}^H$

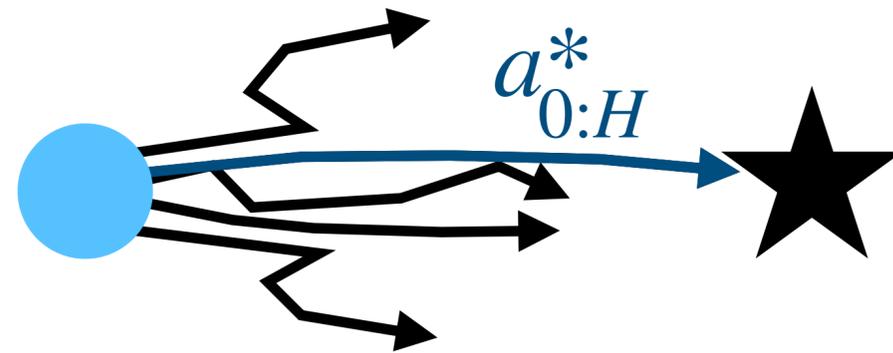
For each iteration

**Sample**  $N$  action sequences  $\{a_{0:H}^i\}_{i=1}^N$  from sampling distribution

**Evaluate** objective  $J(a_{0:H}^i) = \sum_{t=0}^H \gamma^t r(s_t, a_t^i)$  for each sample

**Select** top  $K$  performing samples, i.e. highest value  $J(a_{0:H}^i)$

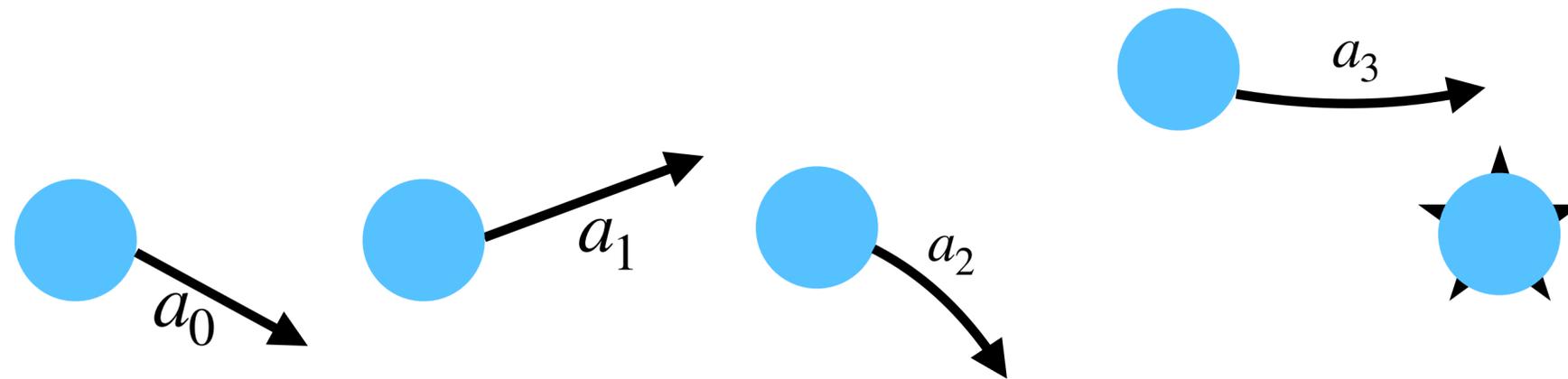
**Update** parameters  $\{\mu_t, \sigma_t^2\}_{t=0}^H$  of action dist. using top  $K$  samples



# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

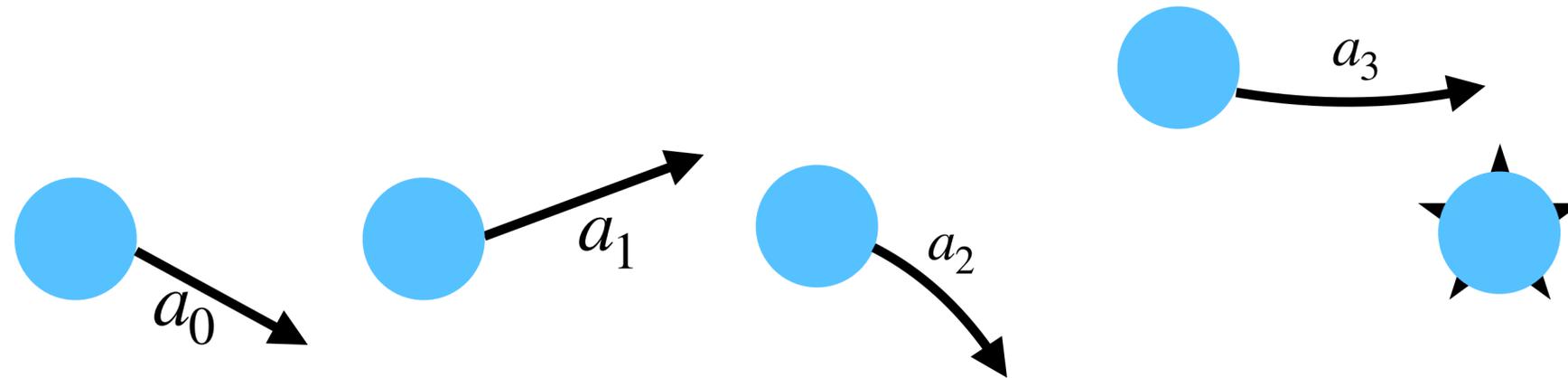


# Collocation methods

## Illustration

$$\boxed{J(a_{0:H}, s_{0:H})} = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

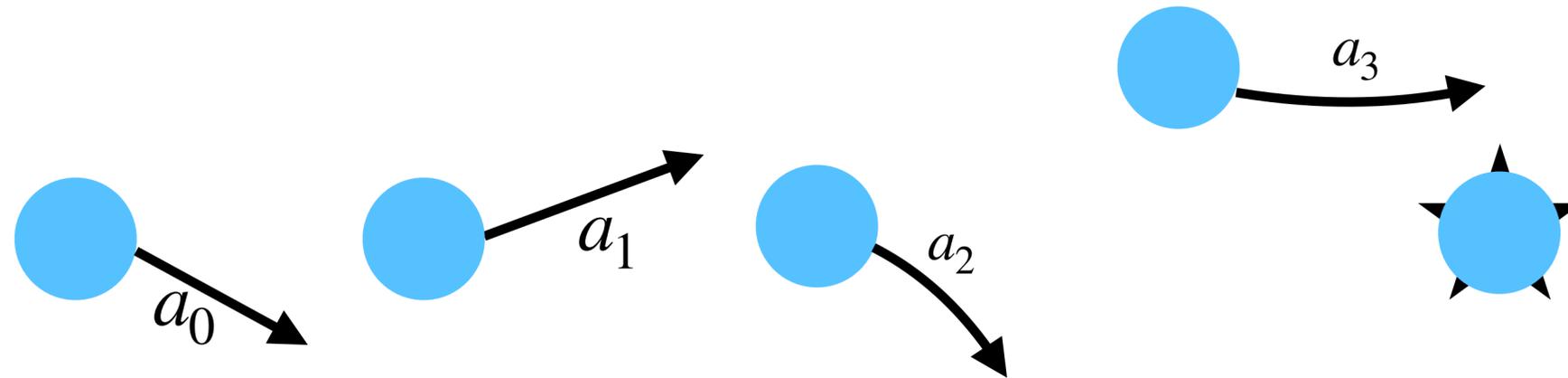


# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions



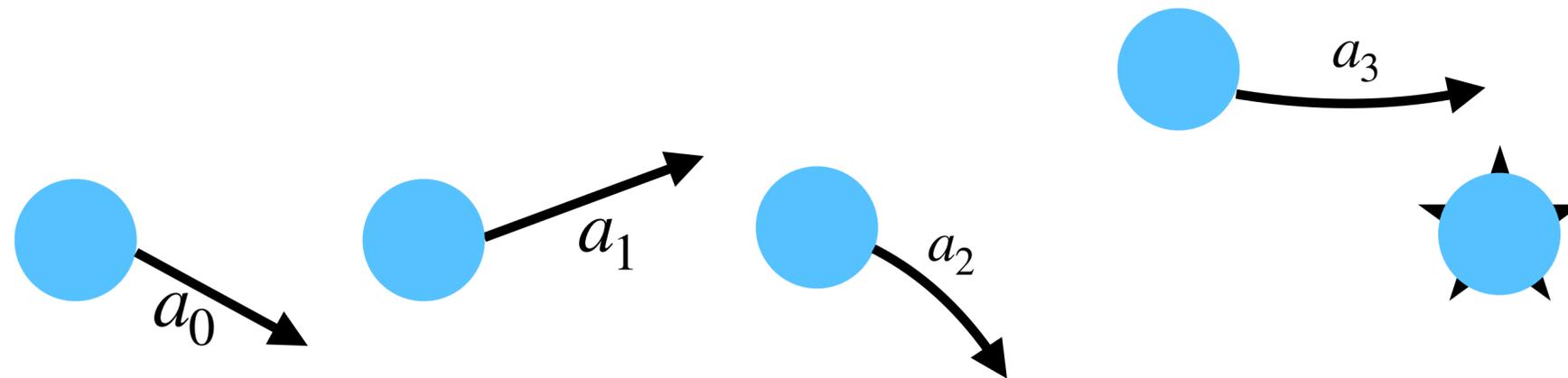
# Collocation methods

## Illustration

$$\boxed{J(a_{0:H}, s_{0:H})} = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \boxed{\|s_{t+1} - f(s_t, a_t)\| = 0}$$

Optimising states and actions

Dynamics constraint



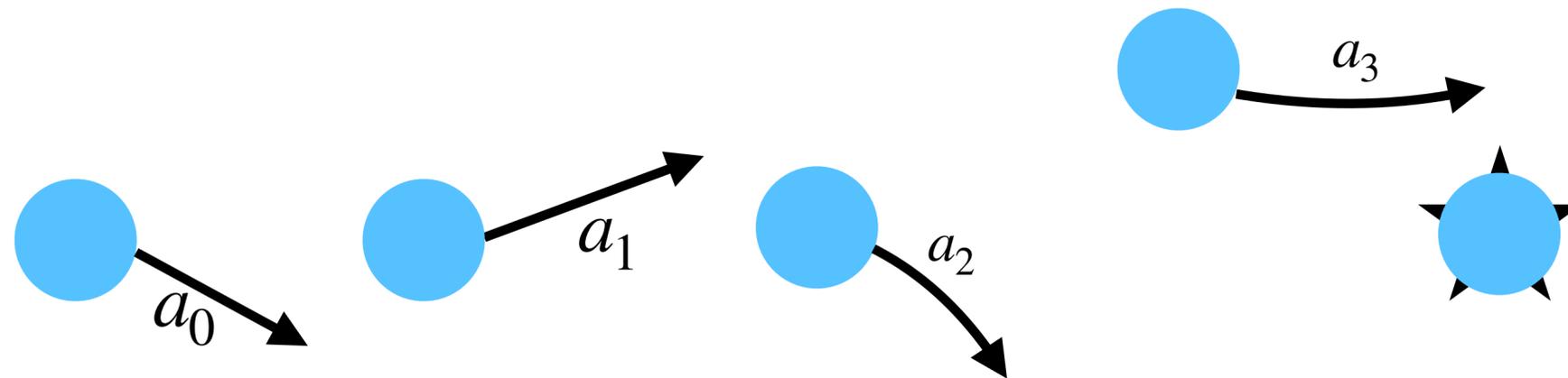
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



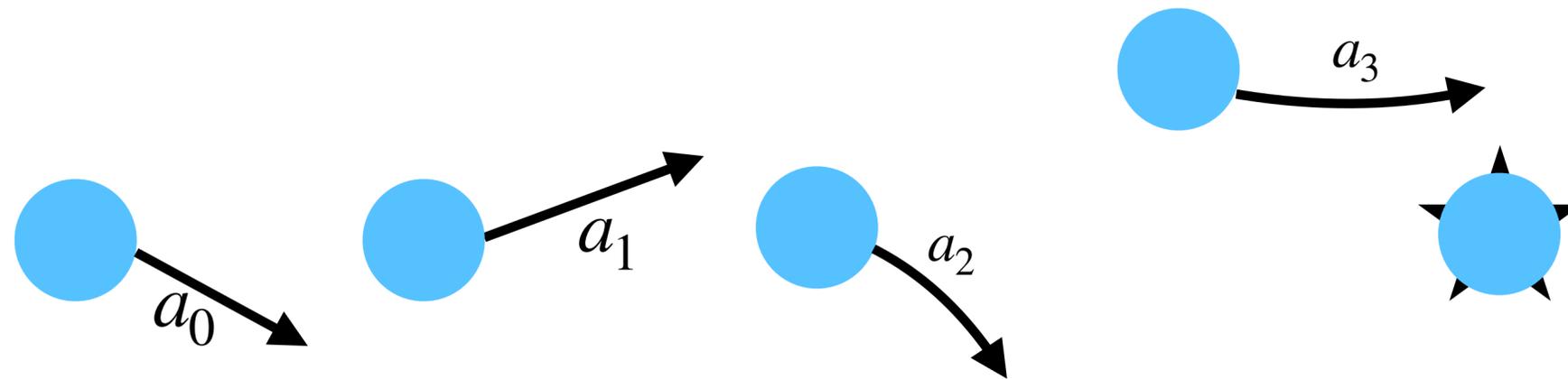
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



Dynamics constraint not satisfied!

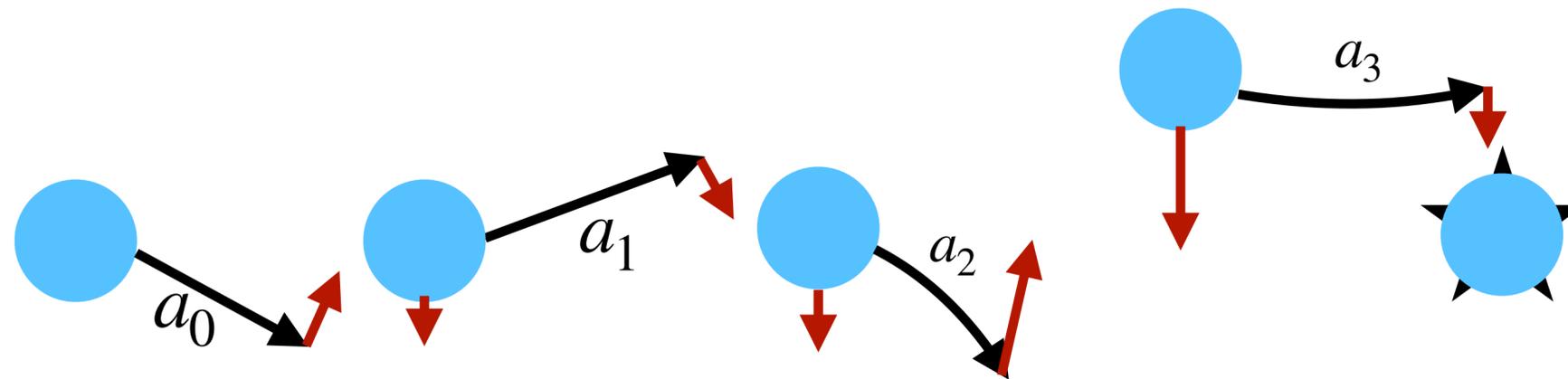
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



Dynamics constraint not satisfied!

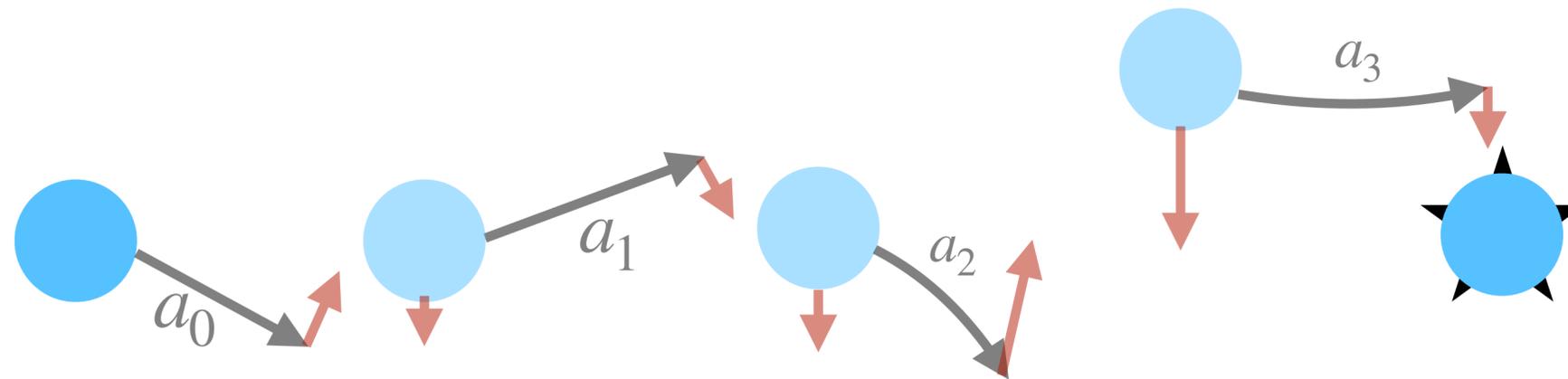
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



Dynamics constraint not satisfied!

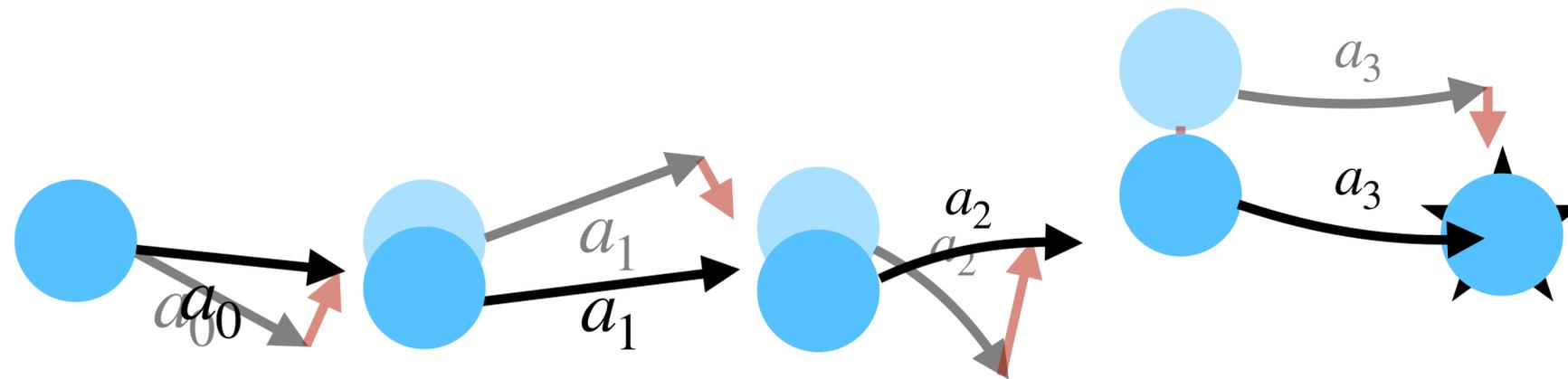
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



Dynamics constraint not satisfied!

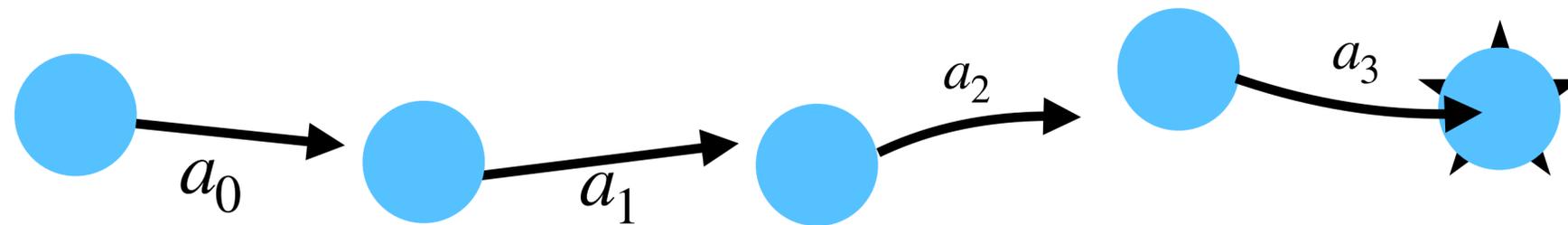
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



Dynamics constraint not satisfied!

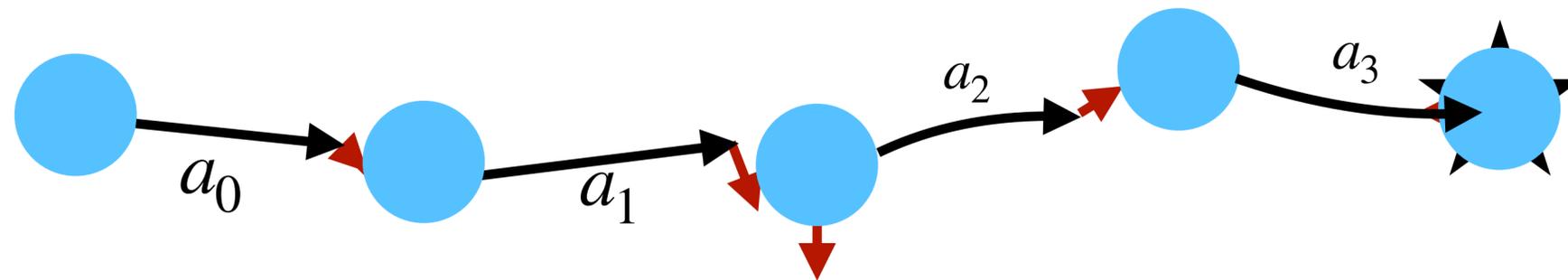
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



Dynamics constraint not satisfied!

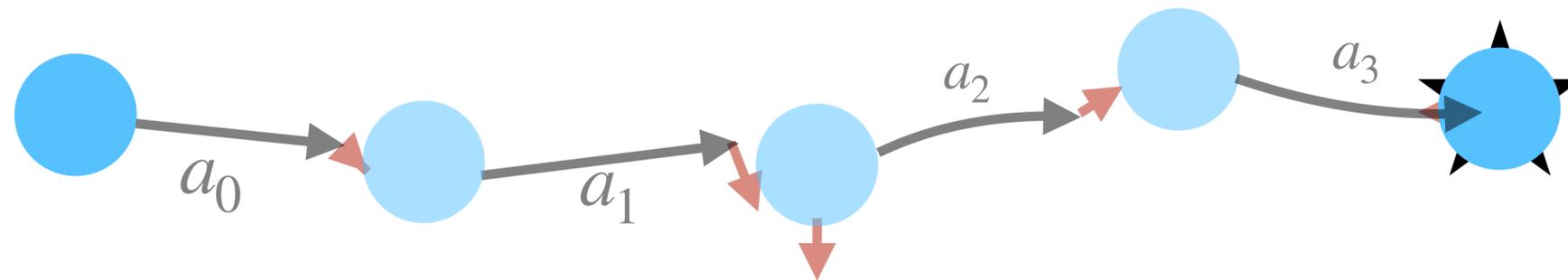
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



Dynamics constraint not satisfied!

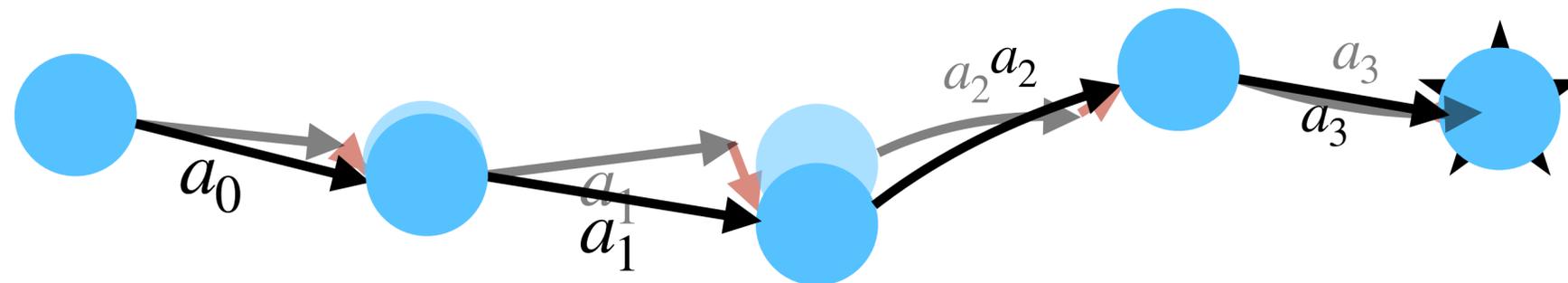
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



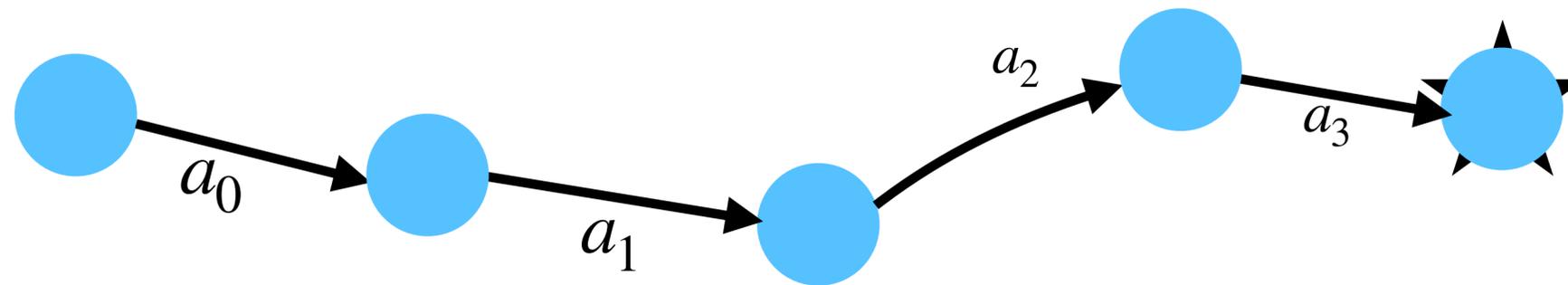
# Collocation methods

## Illustration

$$J(a_{0:H}, s_{0:H}) = \sum_{t=0}^H \gamma^t r(s_t, a_t) \quad \text{s.t.} \quad \|s_{t+1} - f(s_t, a_t)\| = 0$$

Optimising states and actions

Dynamics constraint  
No dynamics rollout



Dynamics constraint satisfied!

# Finite Horizon Planning has Limitations

# Finite Horizon Planning has Limitations

$$\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t)$$

# Finite Horizon Planning has Limitations

$$\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_{\theta}(s_H, a_H)$$

# Finite Horizon Planning has Limitations

$$\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \boxed{\gamma^H Q_{\theta}(s_H, a_H)}$$

Approximate infinite horizon return using learned  $Q$ -function

# Finite Horizon Planning has Limitations

$$\sum_{t=0}^{\infty} \gamma^t Q(s_t, a_t) \approx \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \boxed{\gamma^H Q_{\theta}(s_H, a_H)}$$

Approximate infinite horizon return using learned  $Q$ -function

# Finite Horizon Planning has Limitations

$$\sum_{t=0}^{\infty} \gamma^t Q(s_t, a_t) \approx \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \boxed{\gamma^H Q_{\theta}(s_H, a_H)}$$

Approximate infinite horizon return using learned  $Q$ -function

Learned  $Q$ -function is common in model-free RL

# Finite Horizon Planning has Limitations

$$\sum_{t=0}^{\infty} \gamma^t Q(s_t, a_t) \approx \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \boxed{\gamma^H Q_{\theta}(s_H, a_H)}$$

Approximate infinite horizon return using learned  $Q$ -function

Learned  $Q$ -function is common in model-free RL

**Best of both worlds!**



**Trajectory optimisation methods are open loop.**

**Trajectory optimisation methods are open loop.  
We can do better.**

# Decision-time Planning

## Model Predictive Control (MPC)



# Decision-time Planning

## Model Predictive Control (MPC)

For each environment step



# Decision-time Planning

## Model Predictive Control (MPC)

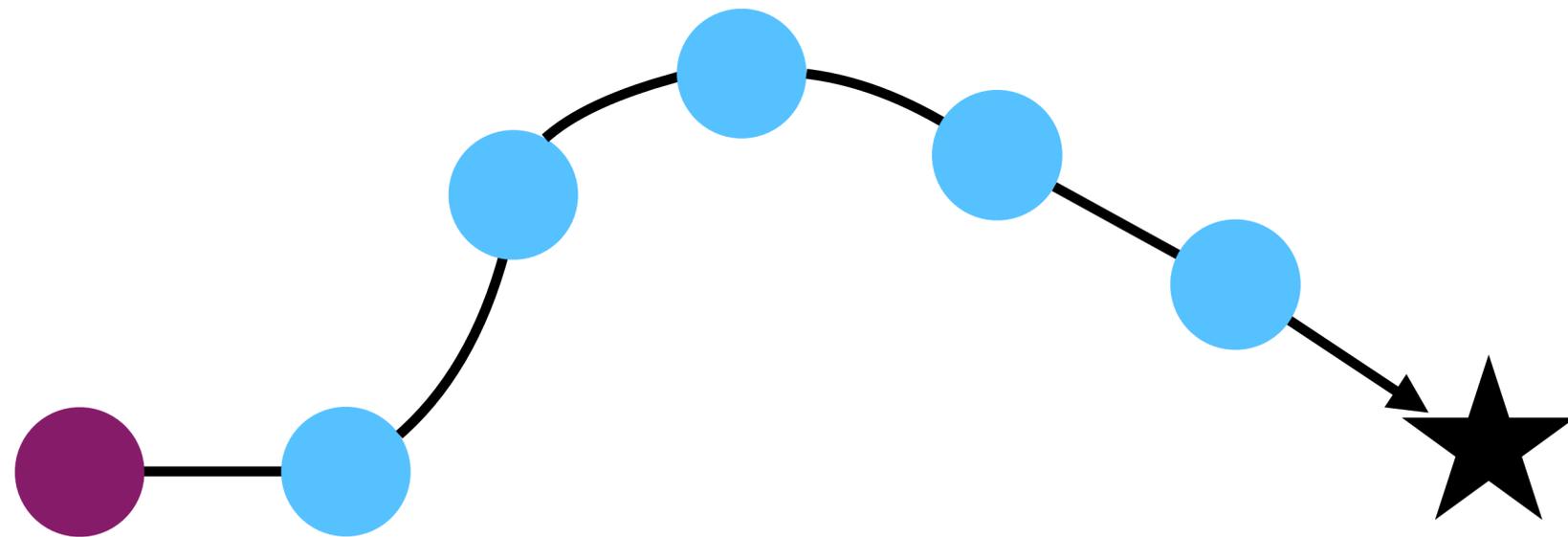
For each environment step

Observe state  $s$



# Decision-time Planning

## Model Predictive Control (MPC)



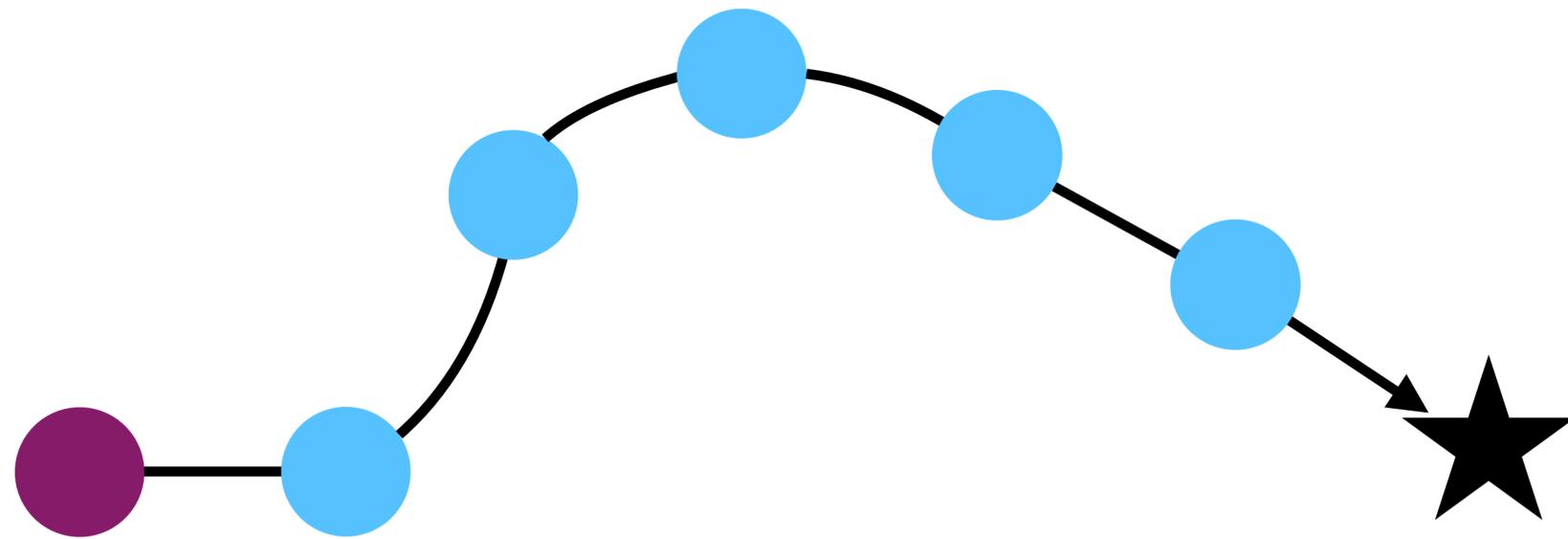
For each environment step

Observe state  $s$

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

# Decision-time Planning

## Model Predictive Control (MPC)



For each environment step

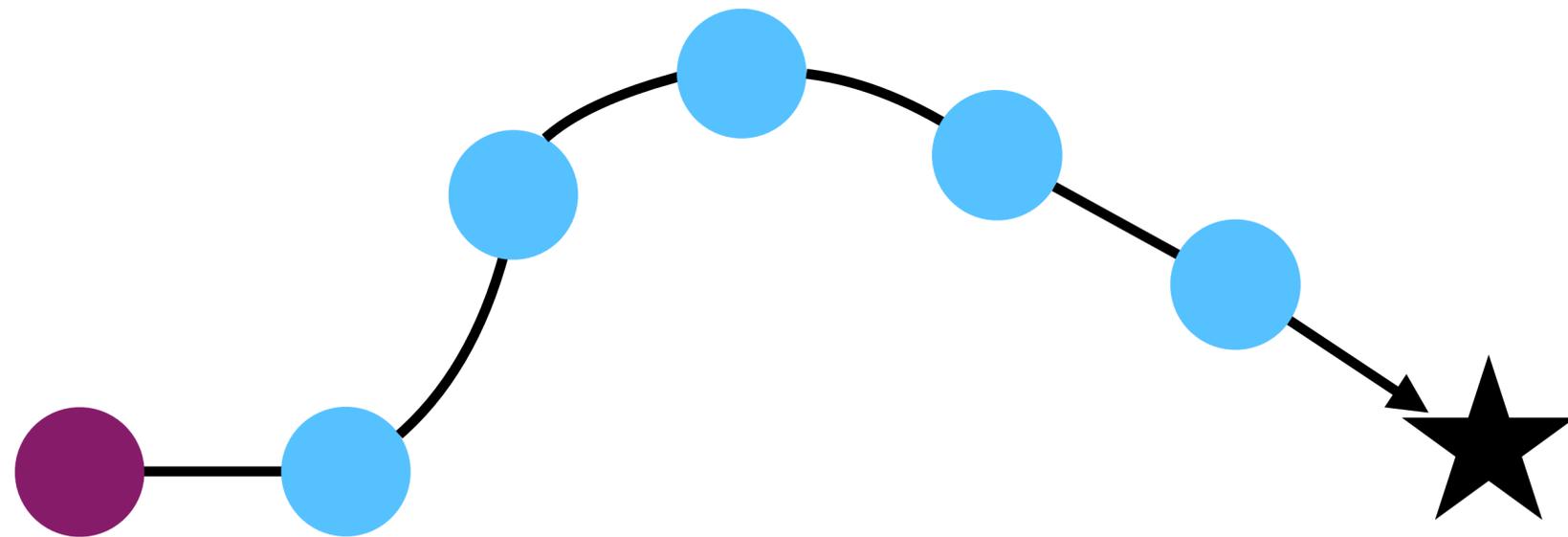
Observe state  $s$

**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

# Decision-time Planning

## Model Predictive Control (MPC)



For each environment step

Observe state  $s$

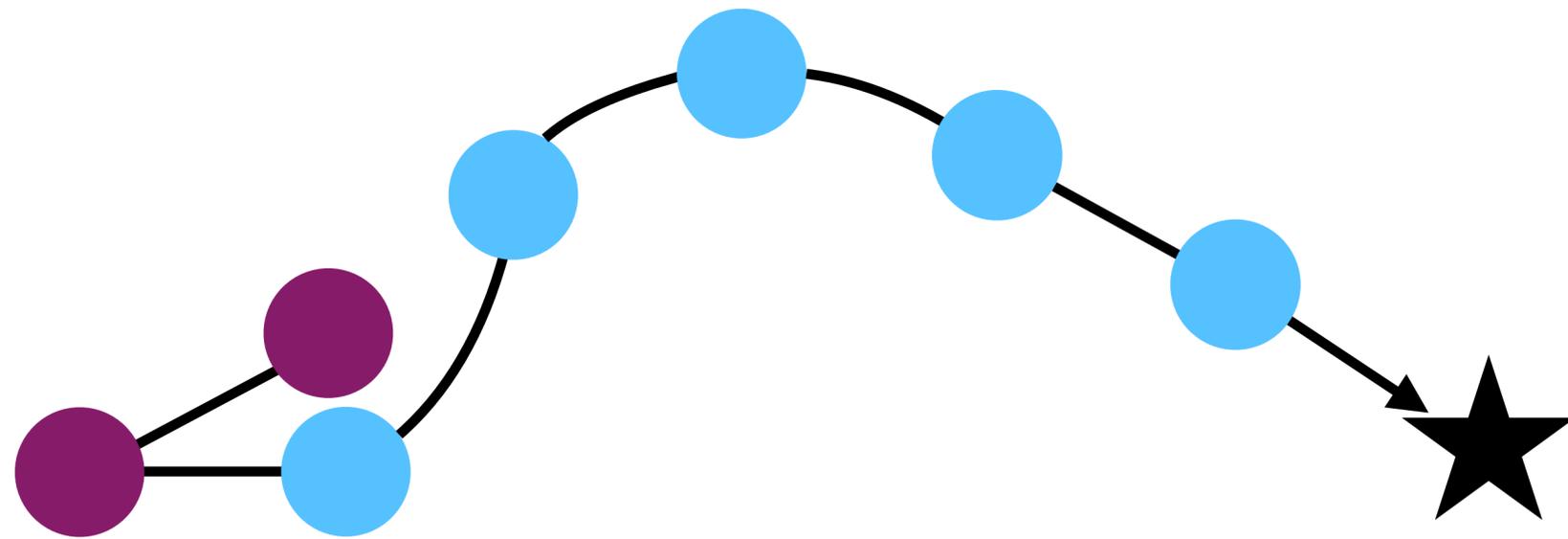
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



For each environment step

Observe state  $s$

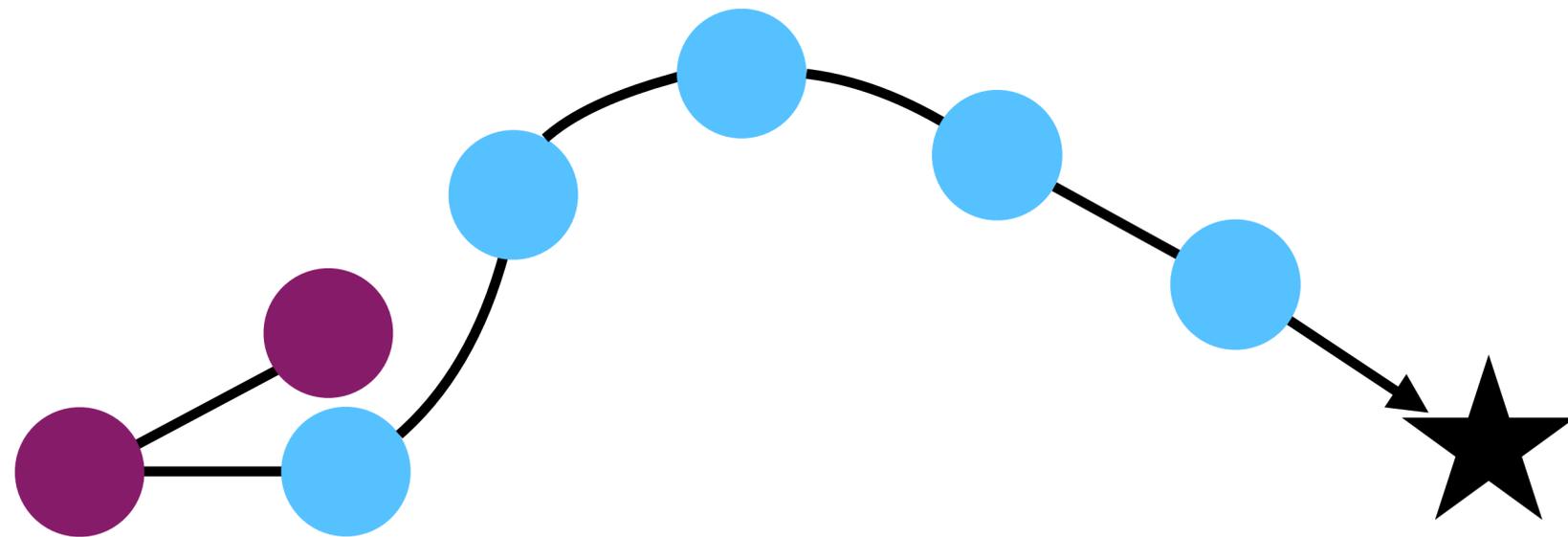
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



For each environment step

Observe state  $s$

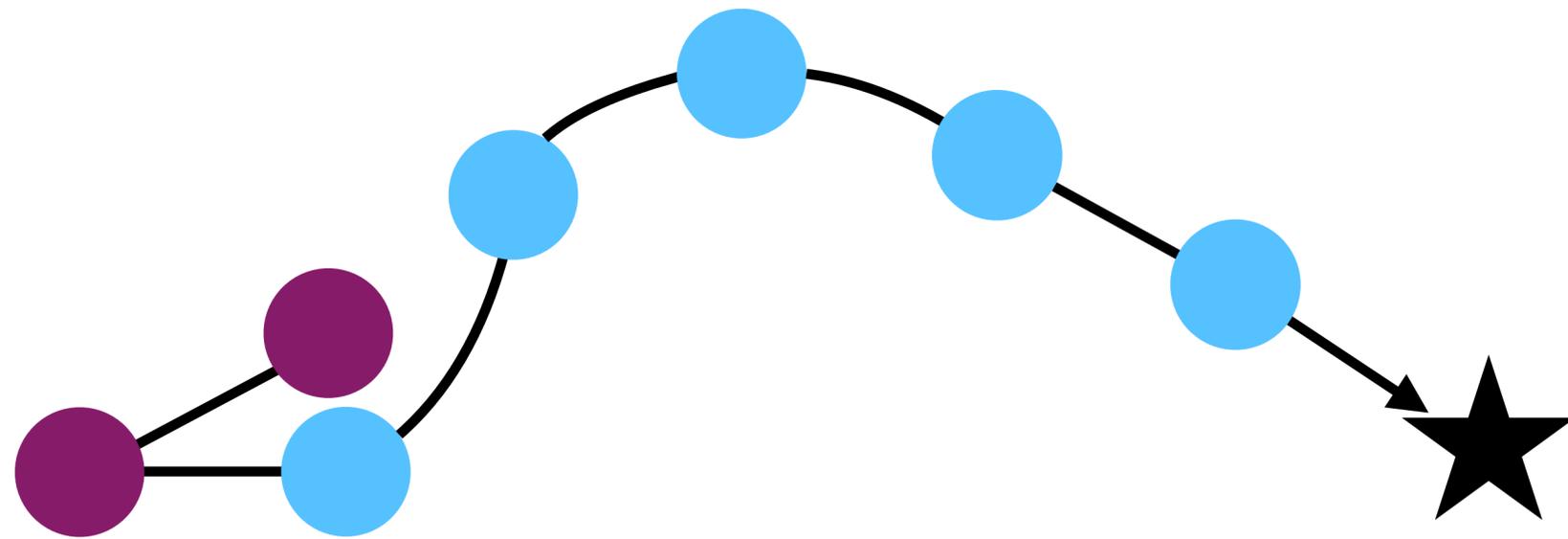
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



**Diverged from planned trajectory...**

**Discard  $a_1, \dots, a_H$**

For each environment step

Observe state  $s$

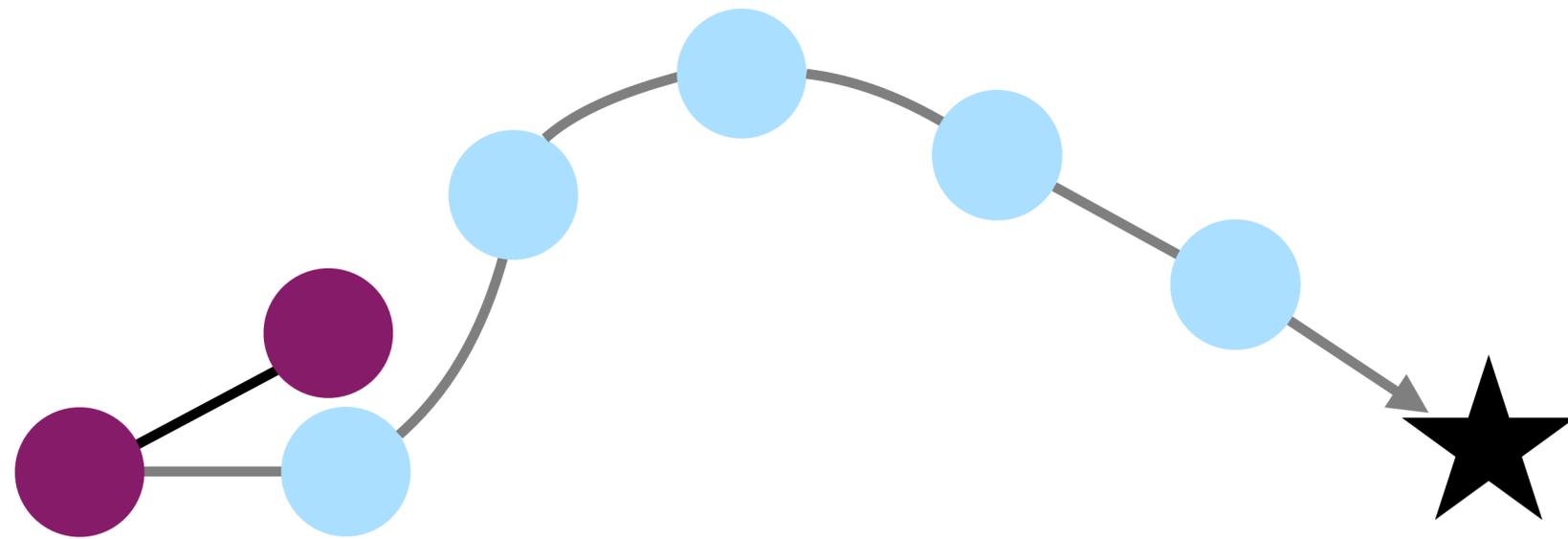
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



**Diverged from planned trajectory...**

**Discard  $a_1, \dots, a_H$**

For each environment step

Observe state  $s$

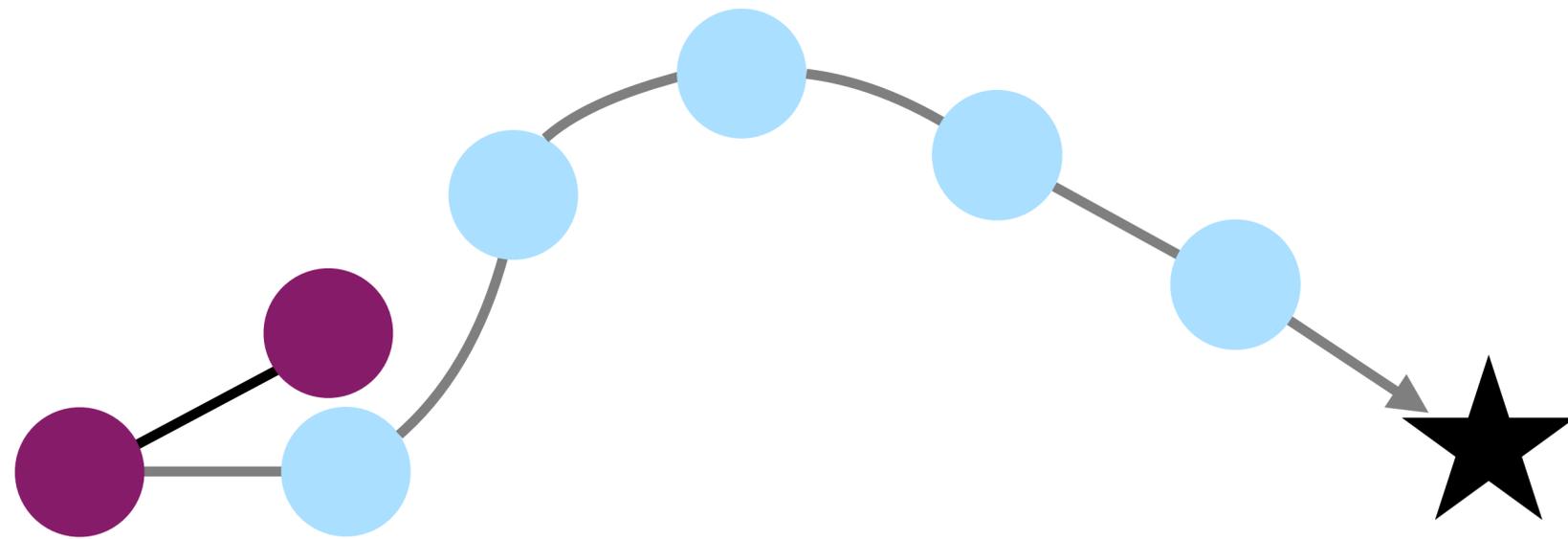
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



**Diverged from planned trajectory...**

**Discard  $a_1, \dots, a_H$**

**So let's replan.**

**FCAI**

For each environment step

Observe state  $s$

**Any trajectory  
optimisation method**

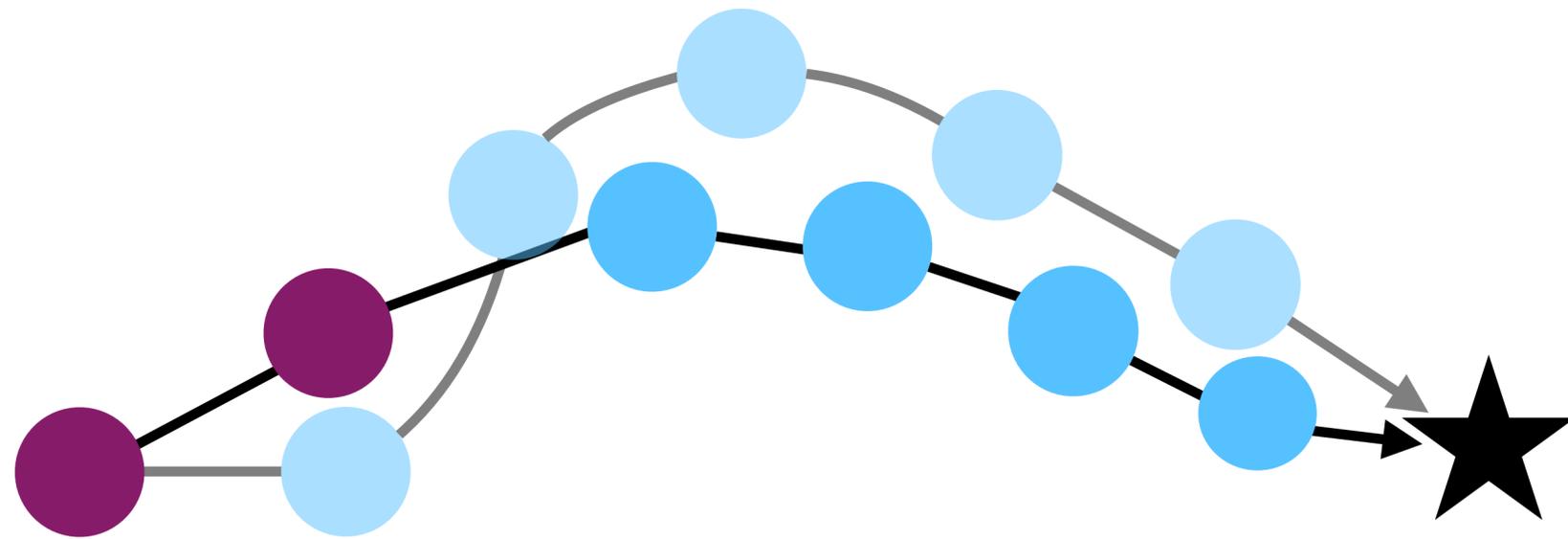
Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

**fcai.fi**

# Decision-time Planning

## Model Predictive Control (MPC)



**Diverged from planned trajectory...**

**Discard  $a_1, \dots, a_H$**

**So let's replan.**

**FCAI**

For each environment step

Observe state  $s$

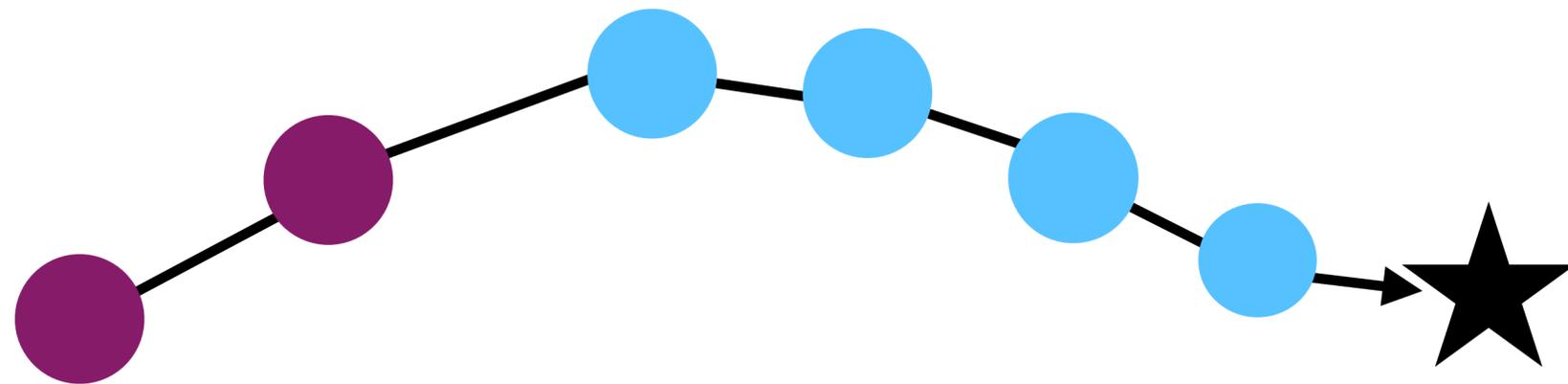
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



Diverged from planned trajectory...

Discard  $a_1, \dots, a_H$

So let's replan.

FCAI

For each environment step

Observe state  $s$

**Any trajectory  
optimisation method**

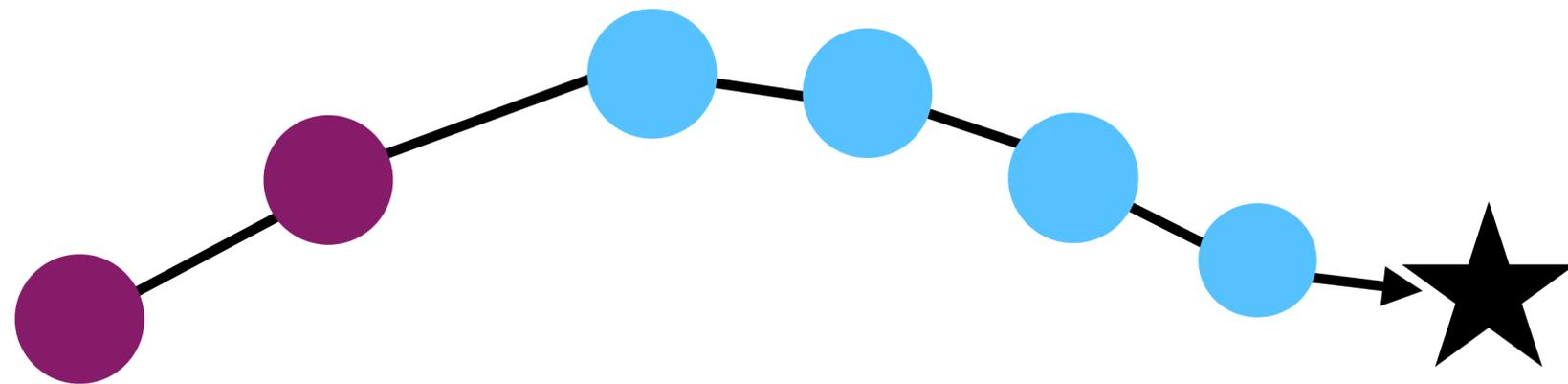
Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

fcai.fi

# Decision-time Planning

## Model Predictive Control (MPC)



For each environment step

Observe state  $s$

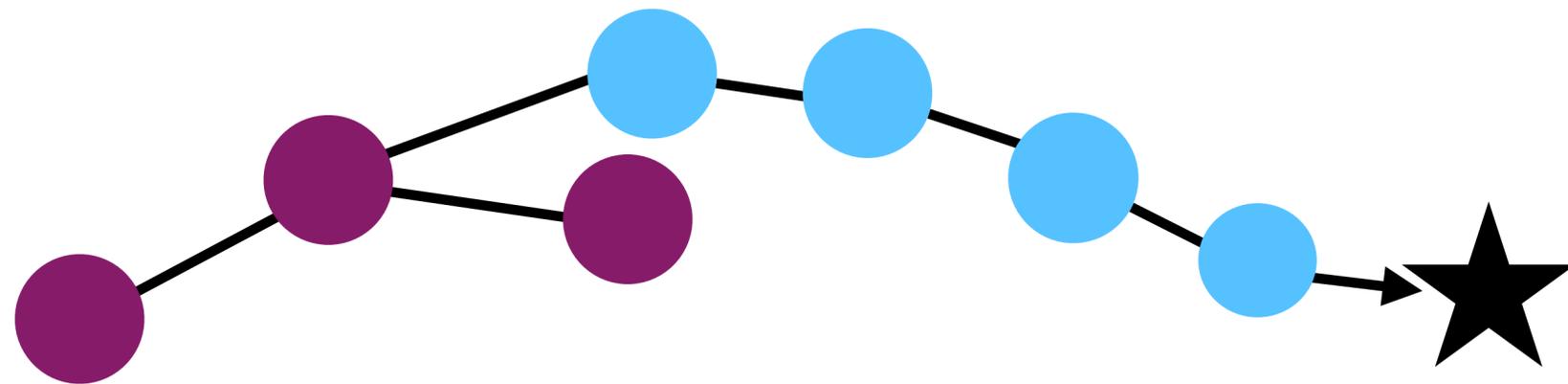
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



For each environment step

Observe state  $s$

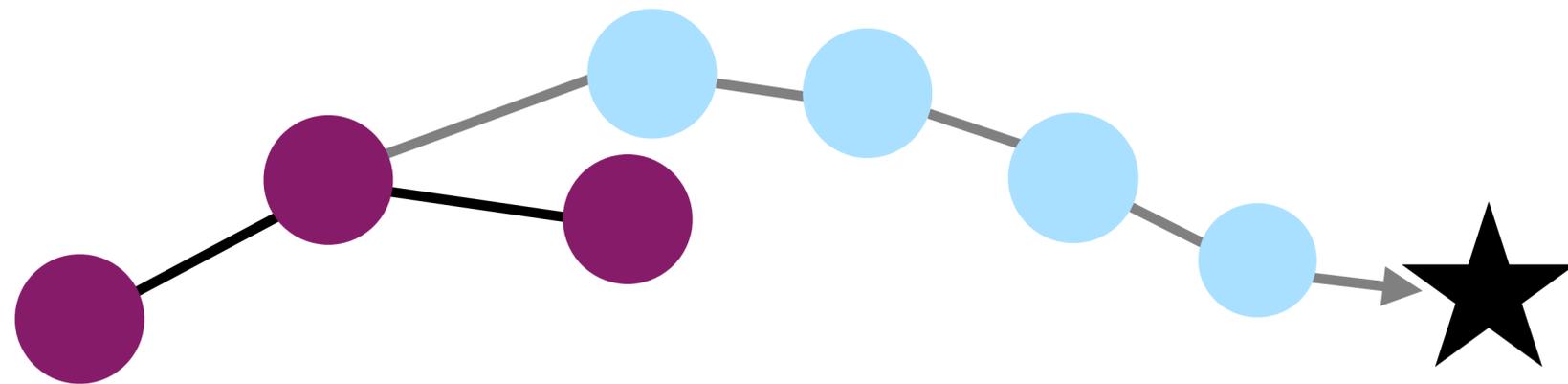
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



For each environment step

Observe state  $s$

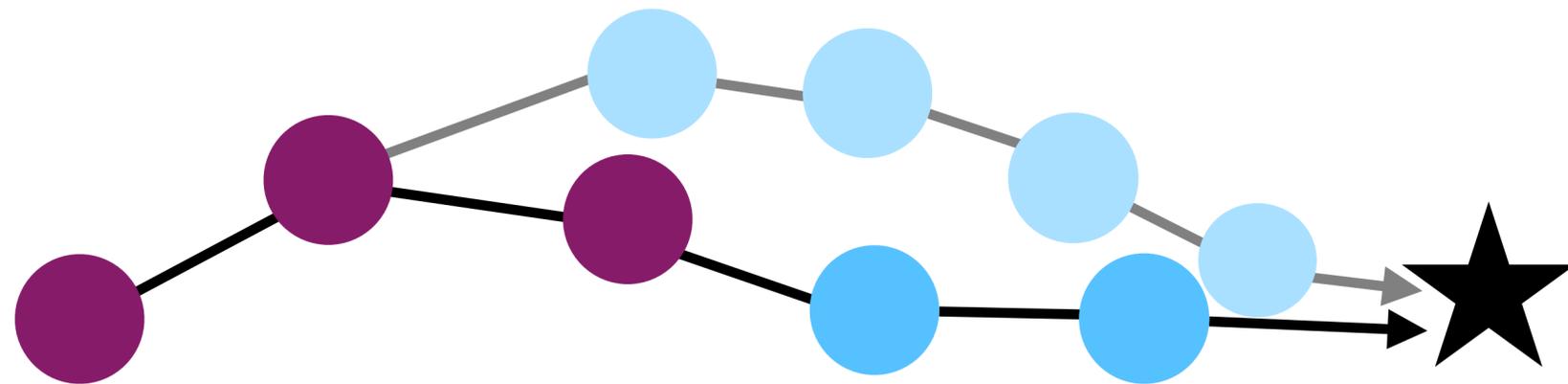
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



For each environment step

Observe state  $s$

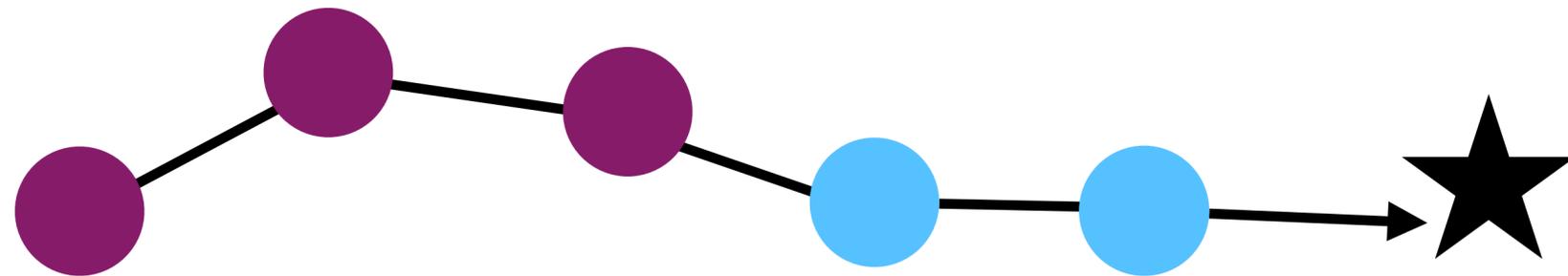
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



For each environment step

Observe state  $s$

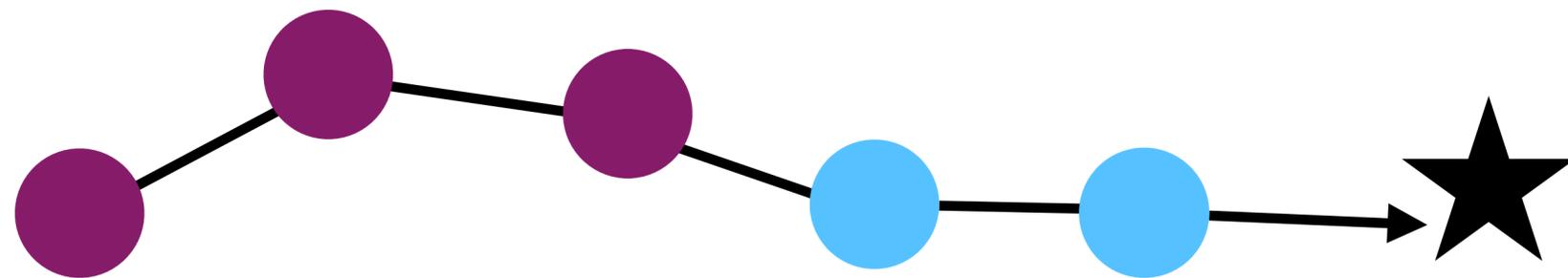
**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)



And so on...

For each environment step

Observe state  $s$

**Any trajectory  
optimisation method**

Plan  $a_{0:H}$  to maximise return  $\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_\theta(s_H, a_H)$

Execute  $a_0$  and discard  $a_1, \dots, a_H$

# Decision-time Planning

## Model Predictive Control (MPC)

$$\pi_{\text{MPC}}(s; f, r, Q_{\theta}) = \arg \max_{a_0} \max_{a_1, \dots, a_{H-1}} = \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H Q_{\theta}(s_H, a_H) \quad \text{s.t.} \quad \begin{aligned} s_{t+1} &= f(s_t, a_t) \\ s_0 &= s \end{aligned}$$

# Decision-time Planning

## Main Takeaways

# Decision-time Planning

## Main Takeaways

Common to use CEM

# Decision-time Planning

## Main Takeaways

Common to use CEM

- Avoids local optima

# Decision-time Planning

## Main Takeaways

Common to use CEM

- Avoids local optima
- Can handle deterministic and stochastic dynamics

# Decision-time Planning

## Main Takeaways

Common to use CEM

- Avoids local optima
- Can handle deterministic and stochastic dynamics
- Avoids exploding/vanishing gradients

# Decision-time Planning

## Main Takeaways

Common to use CEM

- Avoids local optima
- Can handle deterministic and stochastic dynamics
- Avoids exploding/vanishing gradients

Use MPC to make CEM closed loop

# Decision-time Planning

## Main Takeaways

Common to use CEM

- Avoids local optima
- Can handle deterministic and stochastic dynamics
- Avoids exploding/vanishing gradients

Use MPC to make CEM closed loop

Consider infinite horizon via learned  $Q_{\theta}(s, a)$

# Learning Objectives

Understand

1. ~~What a “model” is in model-based RL~~
2. ~~How a “model” can aid decision making~~
3. ~~Differences between background and decision-time planning~~
4. ~~Decision-time planning strategies for continuous actions~~
5. Sources of uncertainty in model-based RL
6. Rationale and insights for decision-making under uncertainty

# Sources of Uncertainty in Model-Based RL

# Sources of Uncertainty

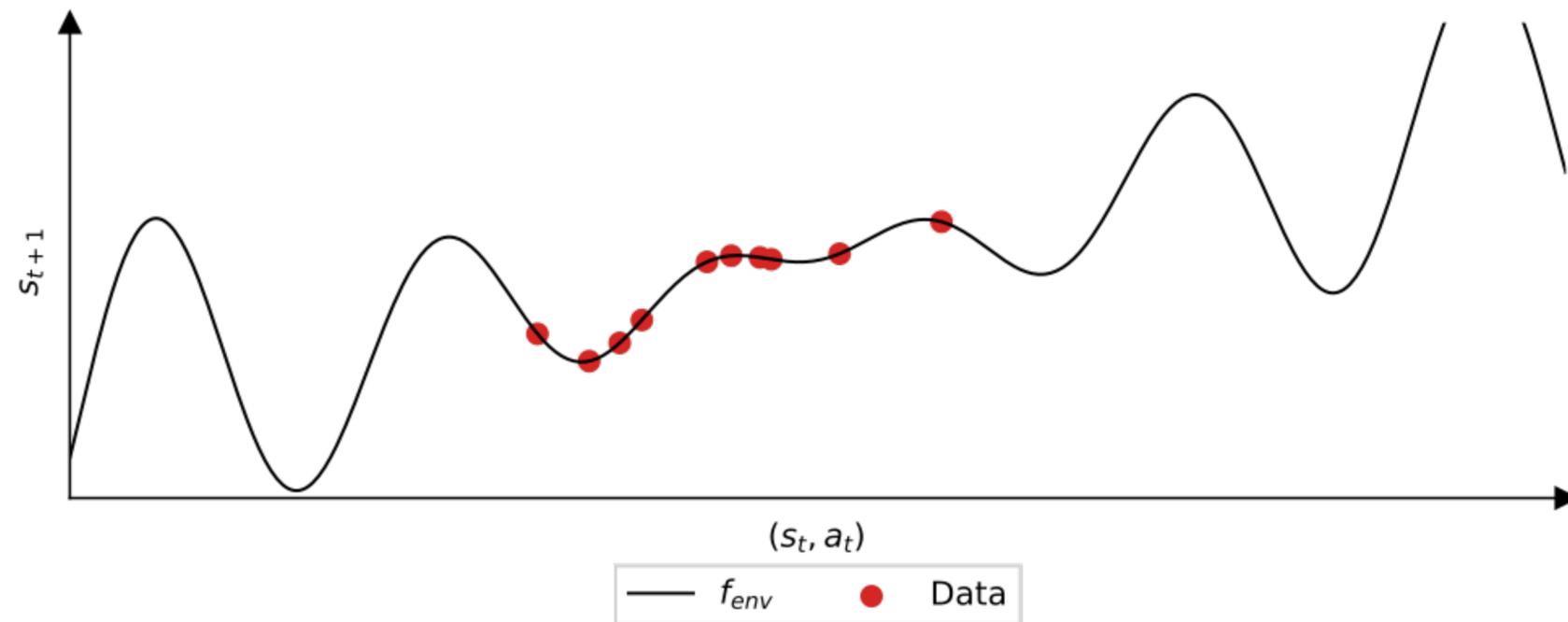
## Learning From Limited Data

$$s_{t+1} = f_{env}(s_t, a_t)$$

# Sources of Uncertainty

## Learning From Limited Data

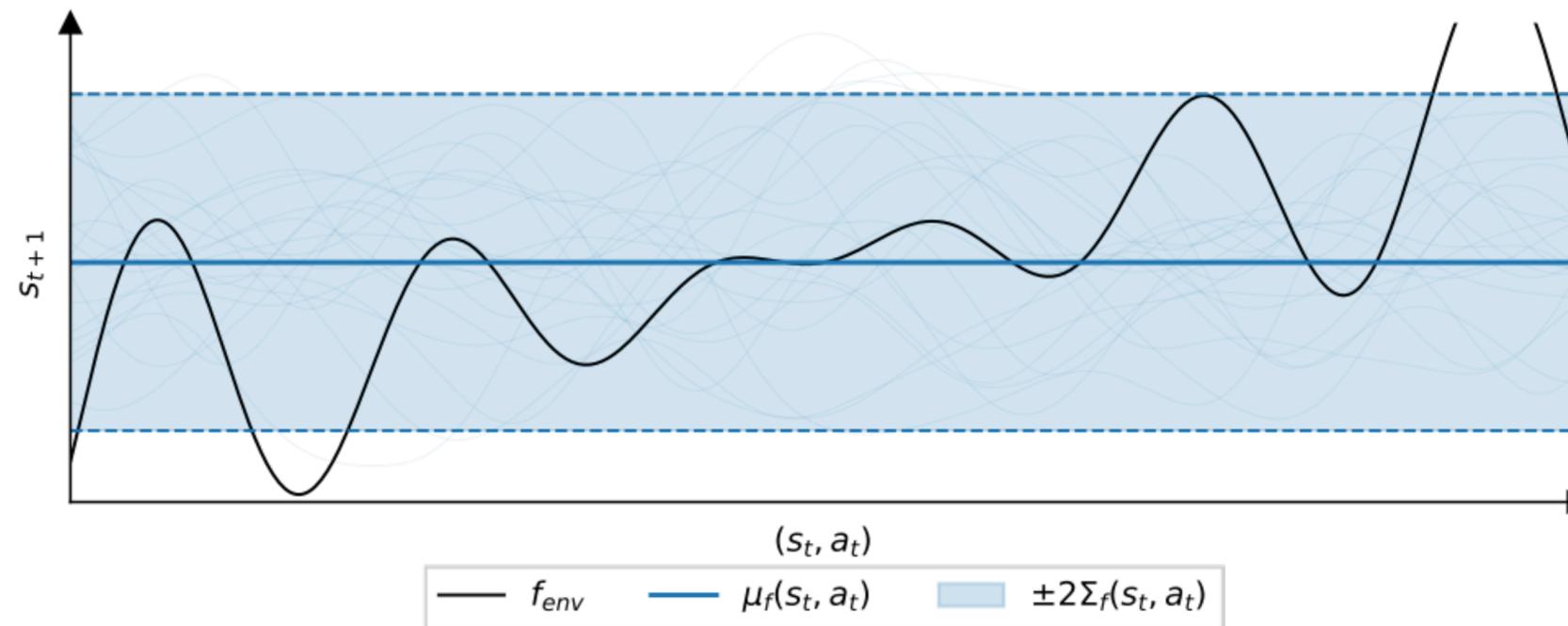
$$s_{t+1} = f_{env}(s_t, a_t)$$



# Sources of Uncertainty

## Learning From Limited Data

$$s_{t+1} = f_{env}(s_t, a_t)$$

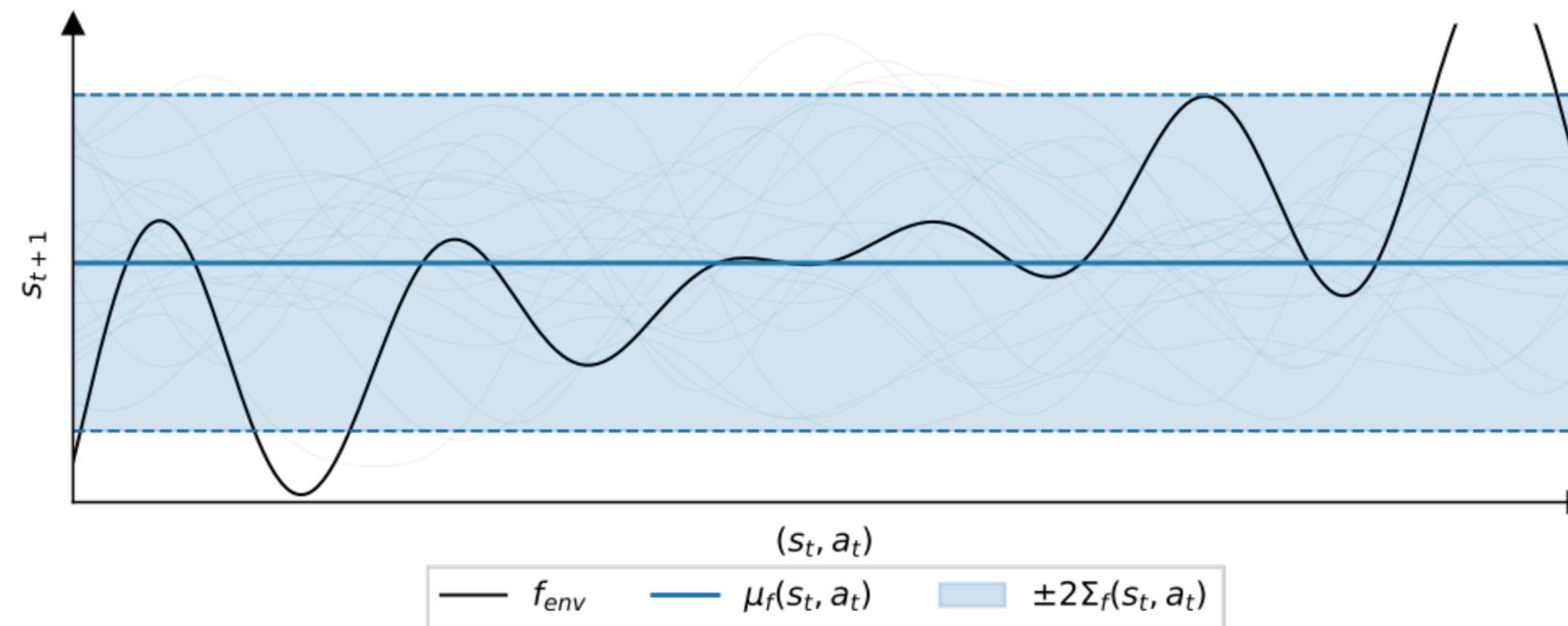


# Sources of Uncertainty

## Learning From Limited Data

$$s_{t+1} = f_{env}(s_t, a_t)$$

Epistemic uncertainty

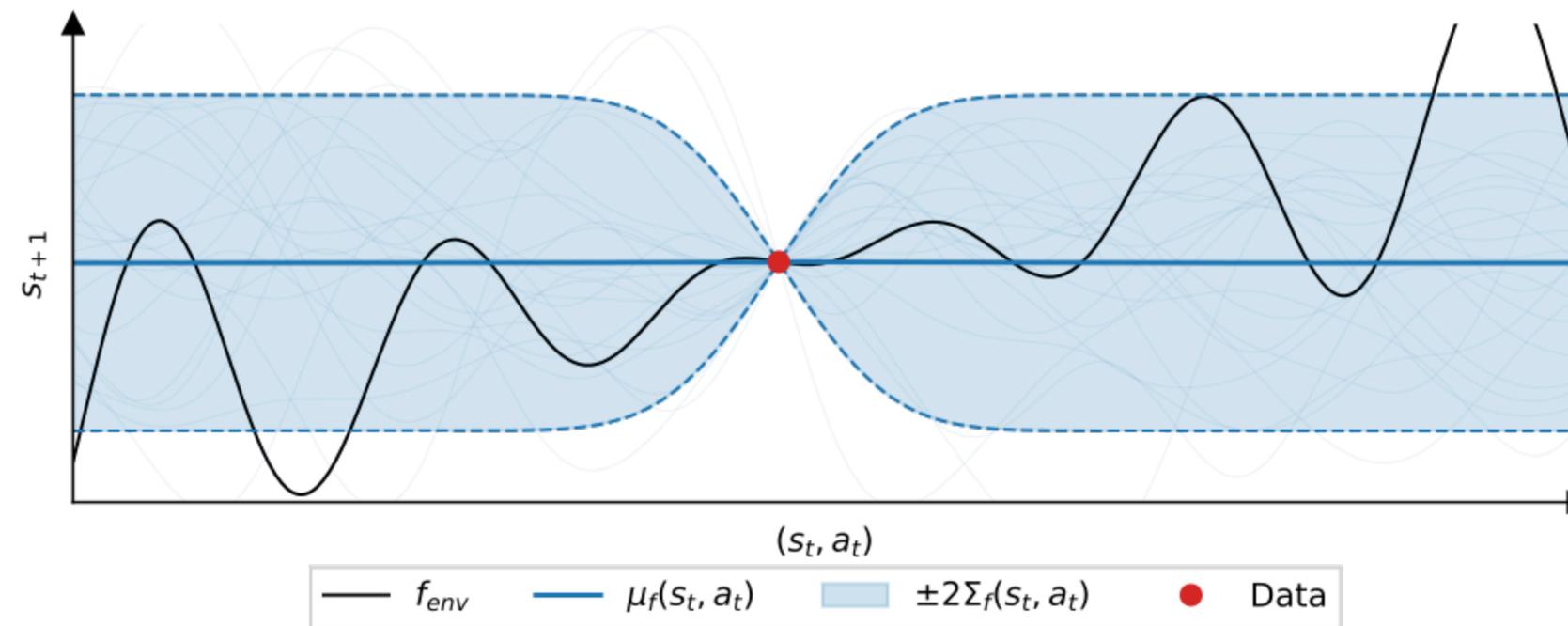


# Sources of Uncertainty

## Learning From Limited Data

$$s_{t+1} = f_{env}(s_t, a_t)$$

Epistemic uncertainty

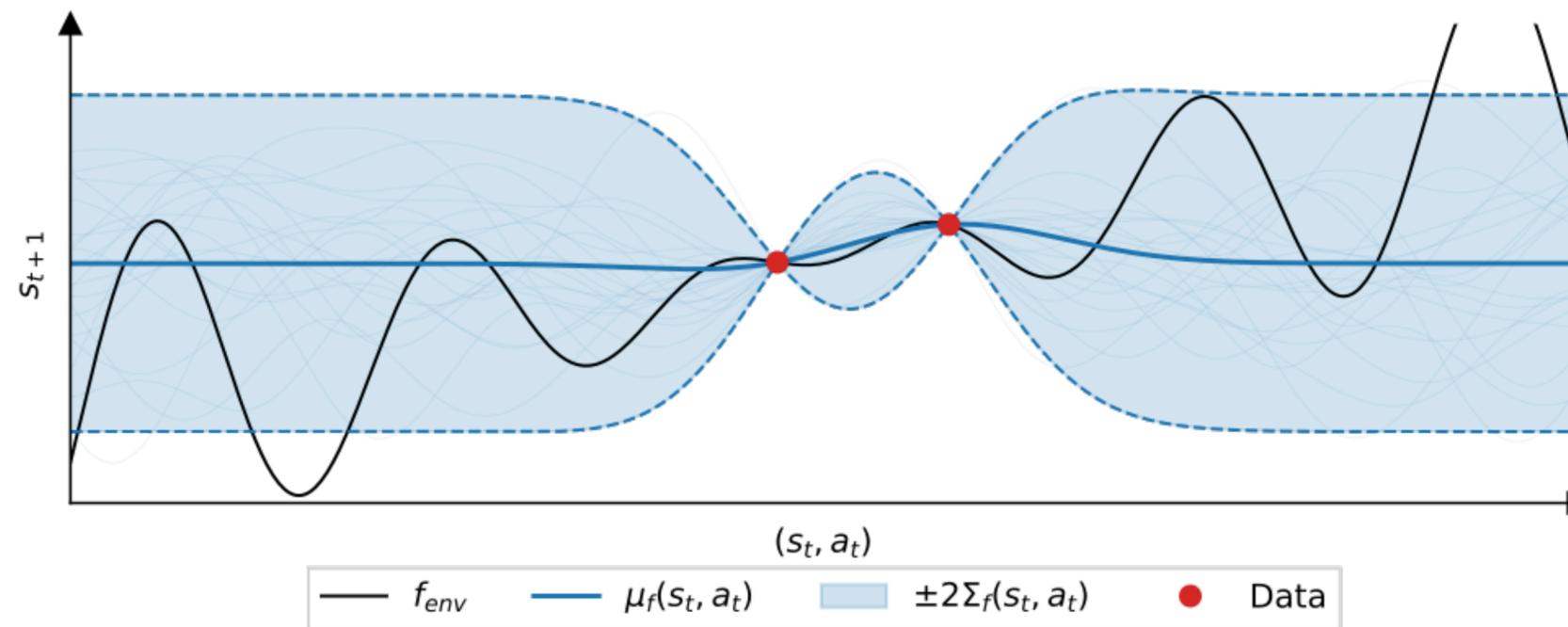


# Sources of Uncertainty

## Learning From Limited Data

$$s_{t+1} = f_{env}(s_t, a_t)$$

Epistemic uncertainty

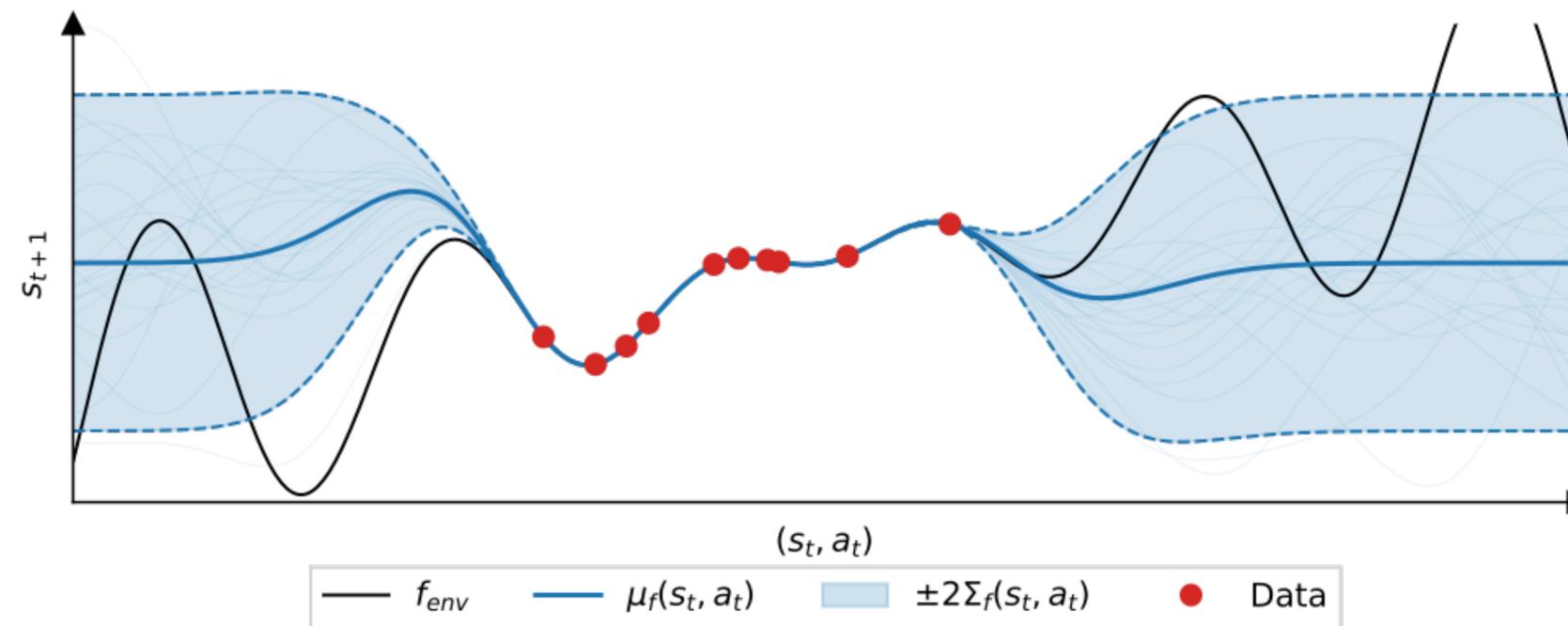


# Sources of Uncertainty

## Learning From Limited Data

$$s_{t+1} = f_{env}(s_t, a_t)$$

Epistemic uncertainty



# Sources of Uncertainty

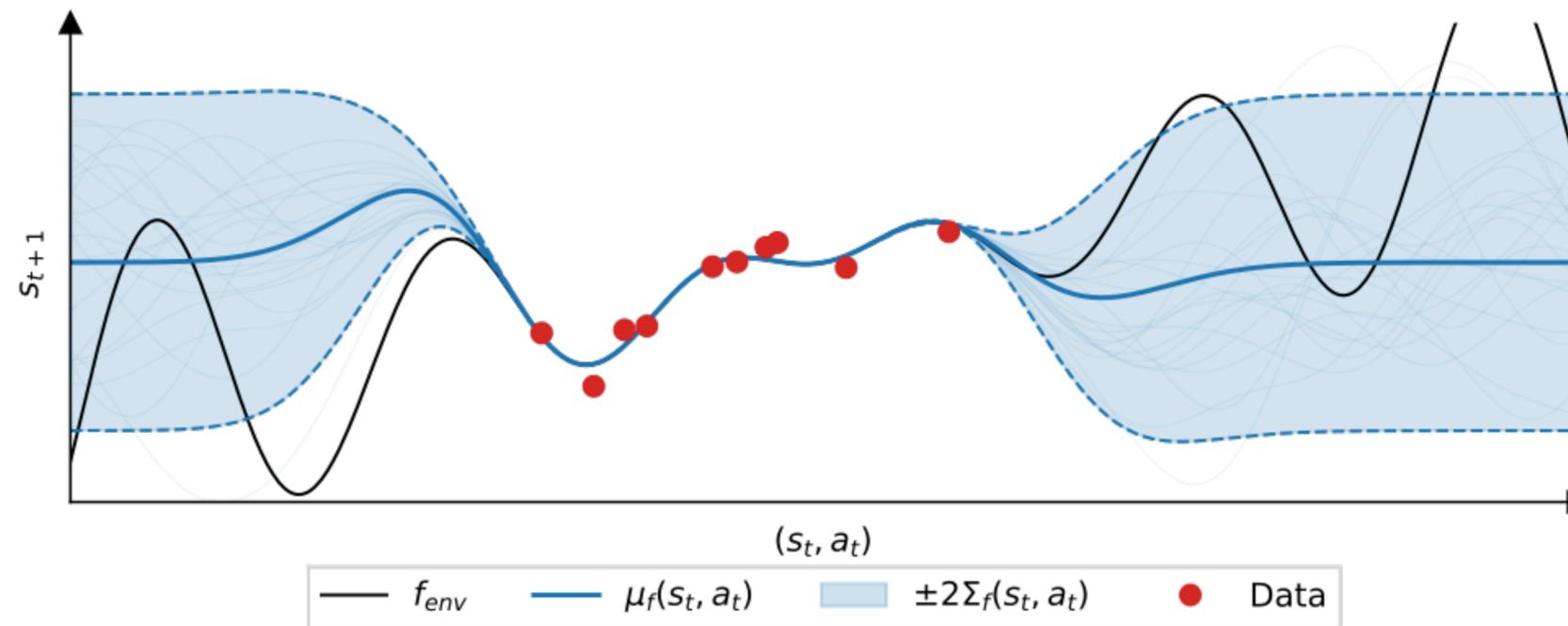
## Stochastic Environments

$$s_{t+1} = f_{env}(s_t, a_t) + \epsilon_t \quad \text{where} \quad \mathbb{E}[\epsilon_t] = 0$$

# Sources of Uncertainty

## Stochastic Environments

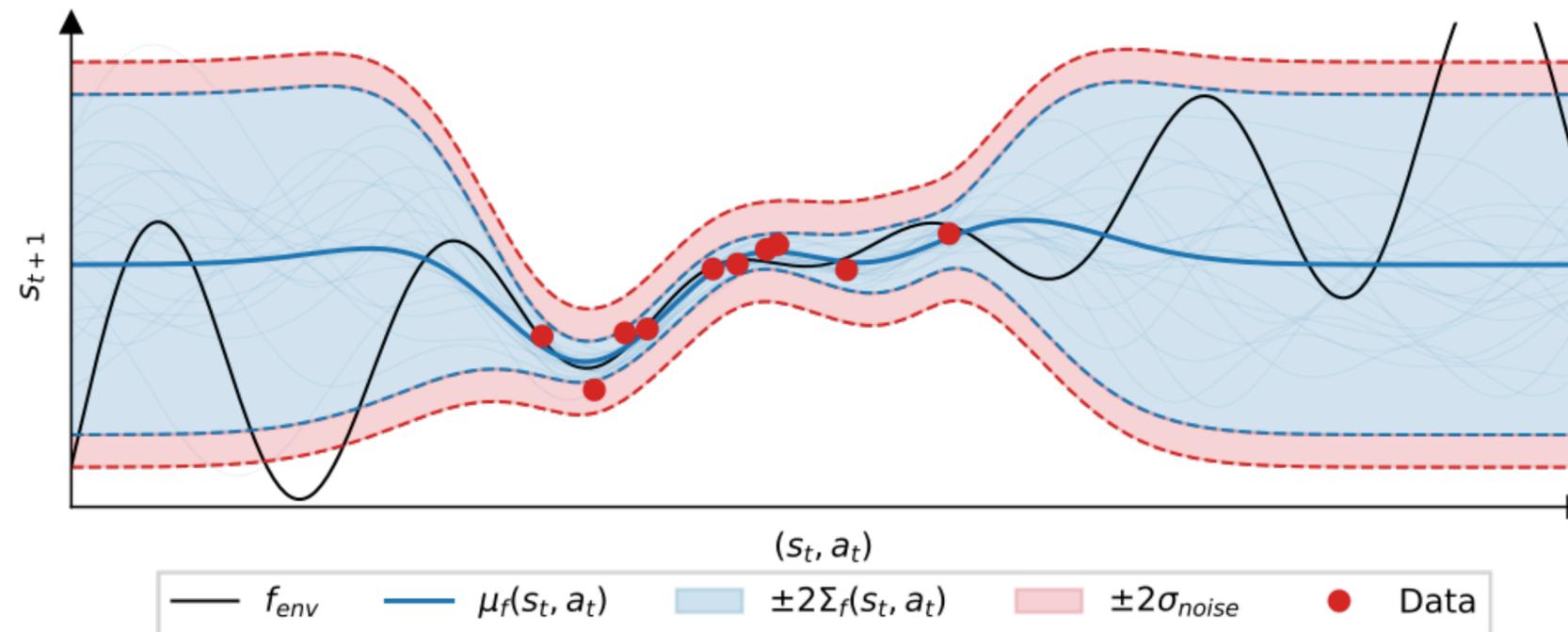
$$s_{t+1} = f_{env}(s_t, a_t) + \epsilon_t \quad \text{where} \quad \mathbb{E}[\epsilon_t] = 0$$



# Sources of Uncertainty

## Stochastic Environments

$$s_{t+1} = f_{env}(s_t, a_t) + \epsilon_t \quad \text{where} \quad \mathbb{E}[\epsilon_t] = 0$$

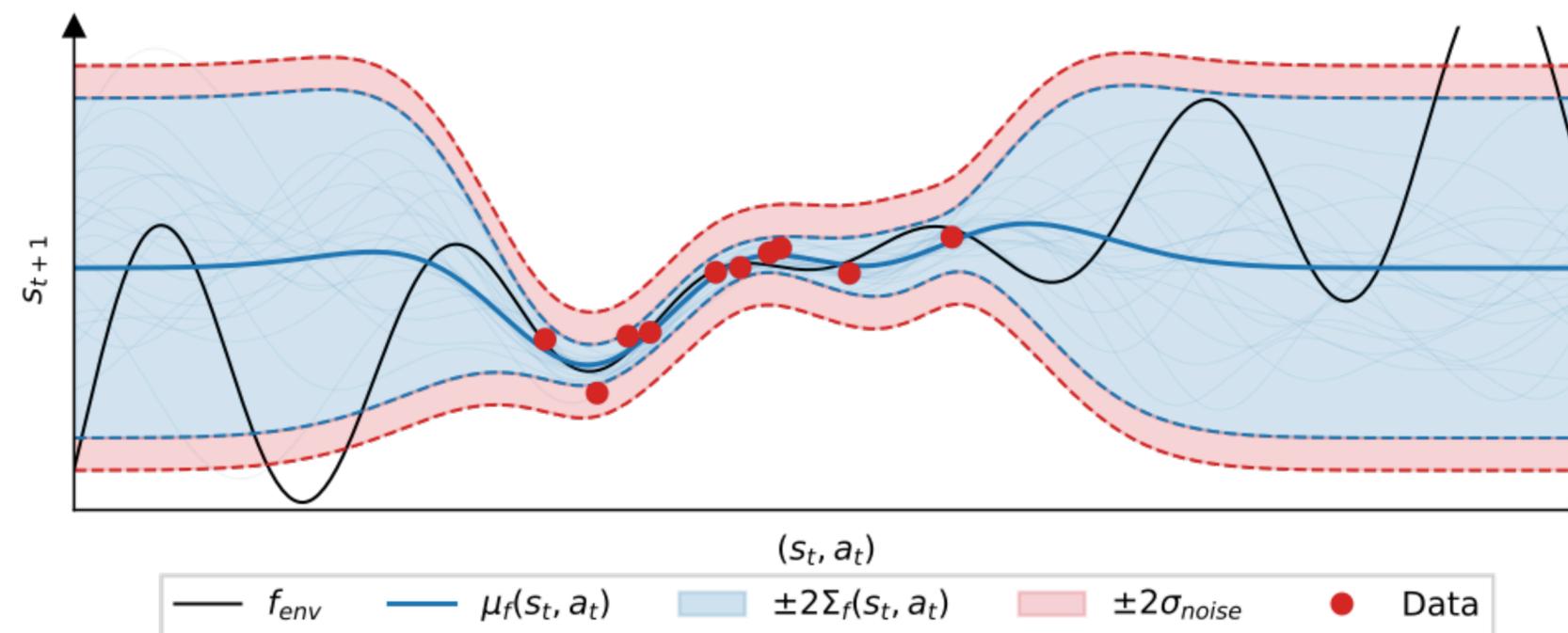


# Sources of Uncertainty

## Stochastic Environments

$$s_{t+1} = f_{env}(s_t, a_t) + \epsilon_t \quad \text{where} \quad \mathbb{E}[\epsilon_t] = 0$$

**Aleatoric uncertainty**



# Decision-making Under Uncertainty

# Sources of Uncertainty

## Decision-making Under Uncertainty

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{???} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, a_t = \pi(s_t) \right]$$

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:** Return = discounted sum of rewards

$$J(\pi; f) = \mathbb{E}_{???} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, a_t = \pi(s_t) \right]$$

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{???} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \overset{\text{Stochastic dynamics}}{s_{t+1} = f(s_t, a_t) + \epsilon_t}, a_t = \pi(s_t) \right]$$

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{???} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, \boxed{a_t = \pi(s_t)} \right]$$

**Deterministic policy**

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{???} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, a_t = \pi(s_t) \right]$$

**What is the expectation over?**

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{???} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, a_t = \pi(s_t) \right]$$

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{\epsilon_{0:\infty}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, a_t = \pi(s_t) \right]$$

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{\epsilon_{0:\infty}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, a_t = \pi(s_t) \right]$$

Expectation is over transition noise, i.e. aleatoric uncertainty

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{\epsilon_{0:\infty}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, a_t = \pi(s_t) \right]$$

Expectation is over transition noise, i.e. aleatoric uncertainty

**Posterior over dynamics models:**

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{\epsilon_{0:\infty}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, a_t = \pi(s_t) \right]$$

Expectation is over transition noise, i.e. aleatoric uncertainty

**Posterior over dynamics models:**

$$p(f \mid \mathcal{D})$$

# Sources of Uncertainty

## Decision-making Under Uncertainty

**RL objective:**

$$J(\pi; f) = \mathbb{E}_{\epsilon_{0:\infty}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_{t+1} = f(s_t, a_t) + \epsilon_t, a_t = \pi(s_t) \right]$$

Expectation is over transition noise, i.e. aleatoric uncertainty

**Posterior over dynamics models:**

$$p(f \mid \mathcal{D})$$

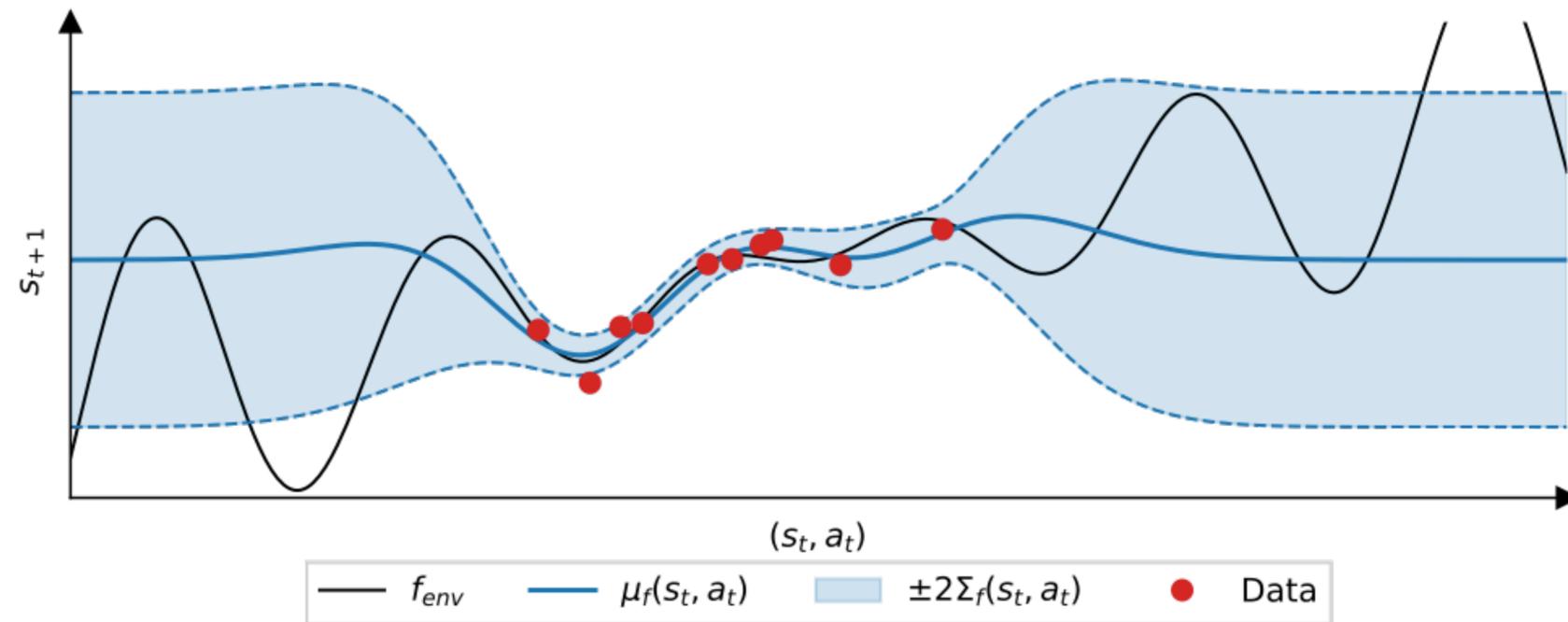
How should we use this?

# Model Averaging

$$\pi^{Greedy} = \arg \max_{\pi} \mathbb{E}_{p(f|\mathcal{D})} [J(\pi; f)]$$

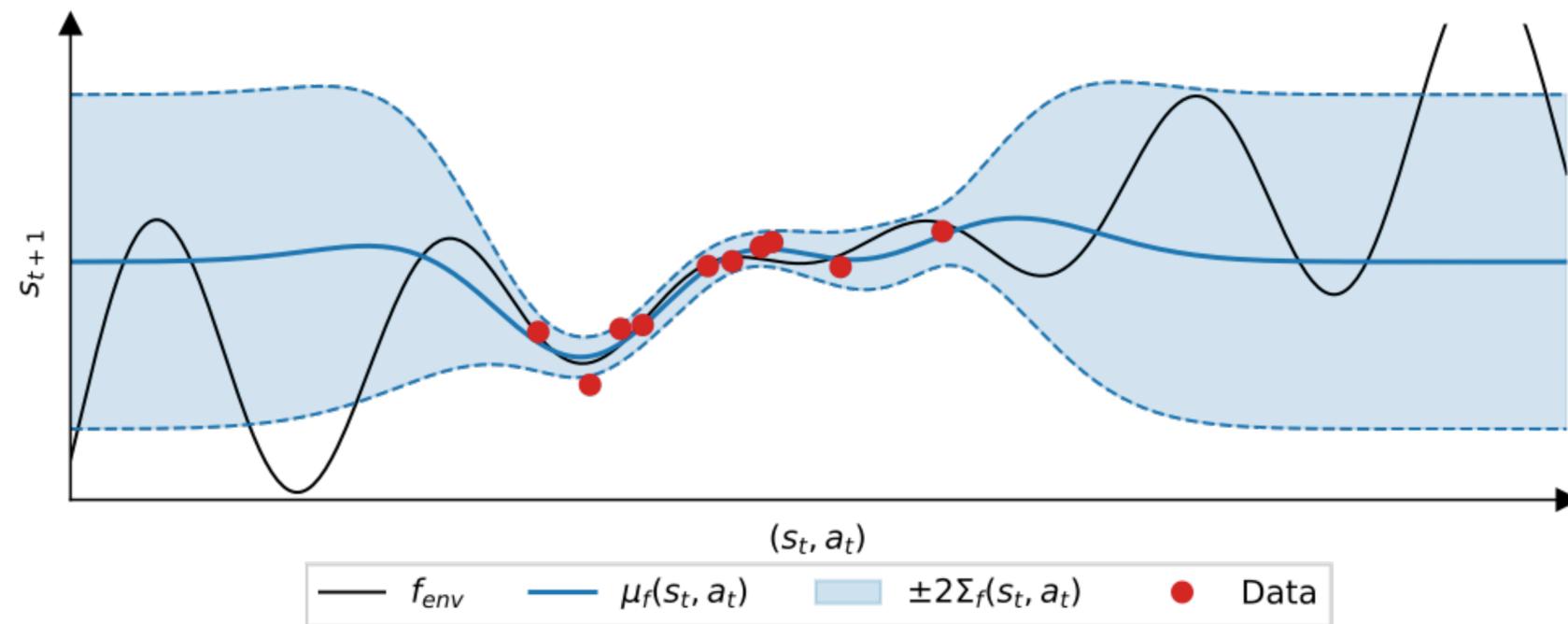
# Model Averaging

$$\pi^{Greedy} = \arg \max_{\pi} \mathbb{E}_{p(f|\mathcal{D})} [J(\pi; f)]$$



# Model Averaging

$$\pi^{Greedy} = \arg \max_{\pi} \mathbb{E}_{p(f|\mathcal{D})} [J(\pi; f)]$$



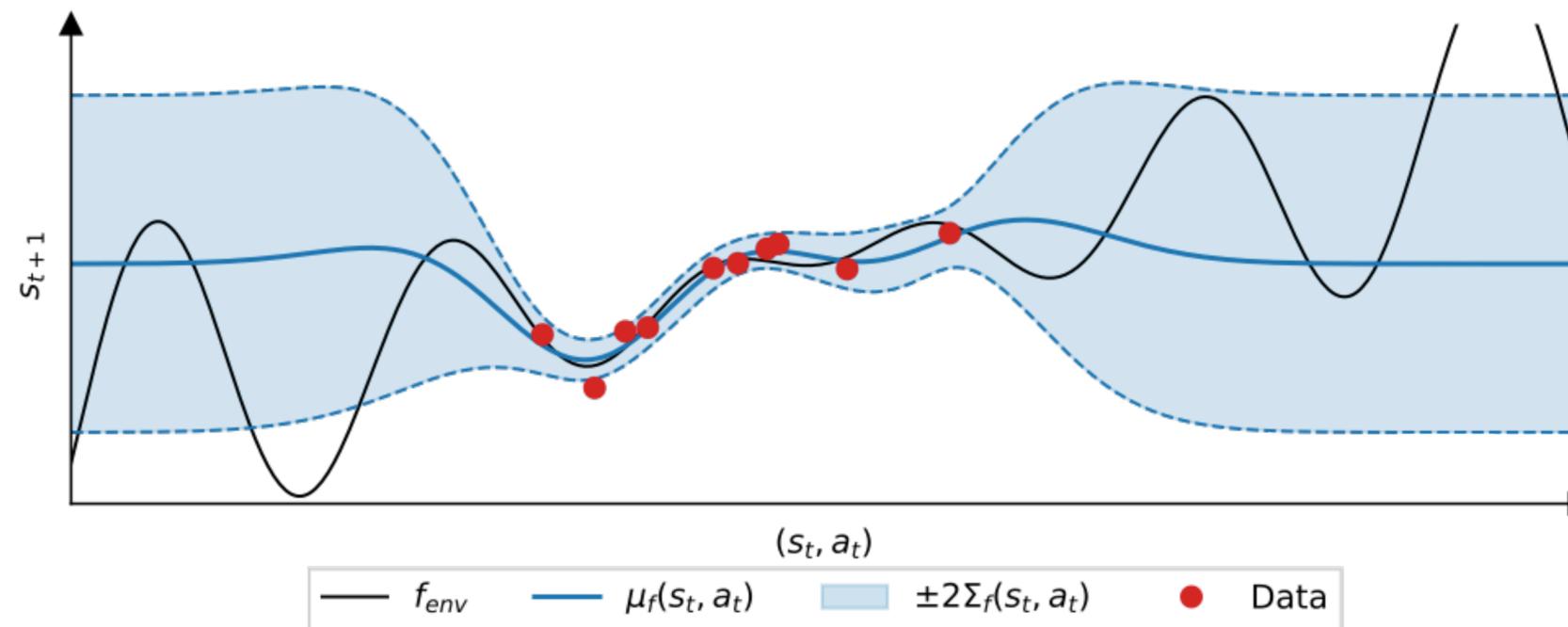
Combats model bias

# Model Averaging

$$\pi^{Greedy} = \arg \max_{\pi} \mathbb{E}_{p(f|\mathcal{D})} [J(\pi; f)]$$

PILCO, PETS, etc

Combats model bias

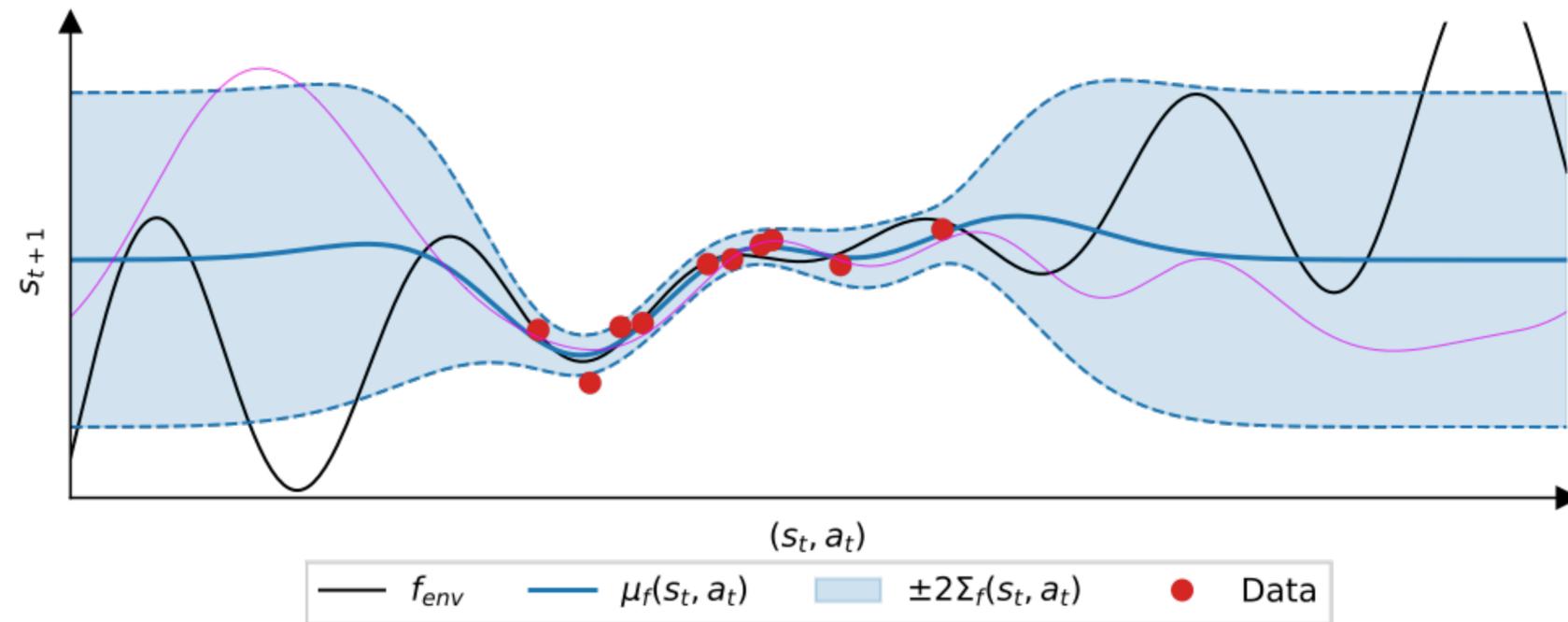


Deisenroth et al. (2011). PILCO: A Model-Based and Data-Efficient Approach to Policy Search. ICML.

Kurtland et al. (2018). Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models. NeurIPS. [fcai.fi](http://fcai.fi)

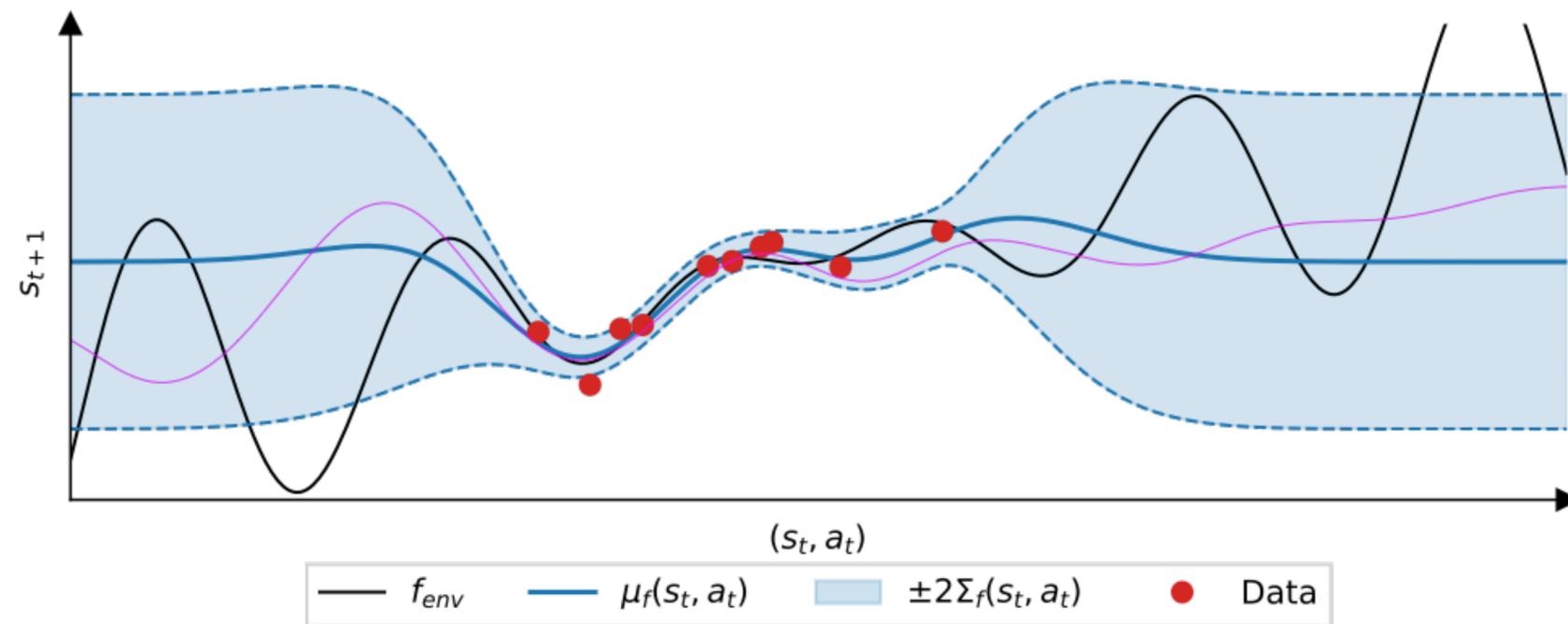
# Exploration via Posterior Sampling

$$\pi^{PS} = \arg \max_{\pi} J(\pi; \tilde{f}), \quad \tilde{f} \sim p(f | \mathcal{D})$$



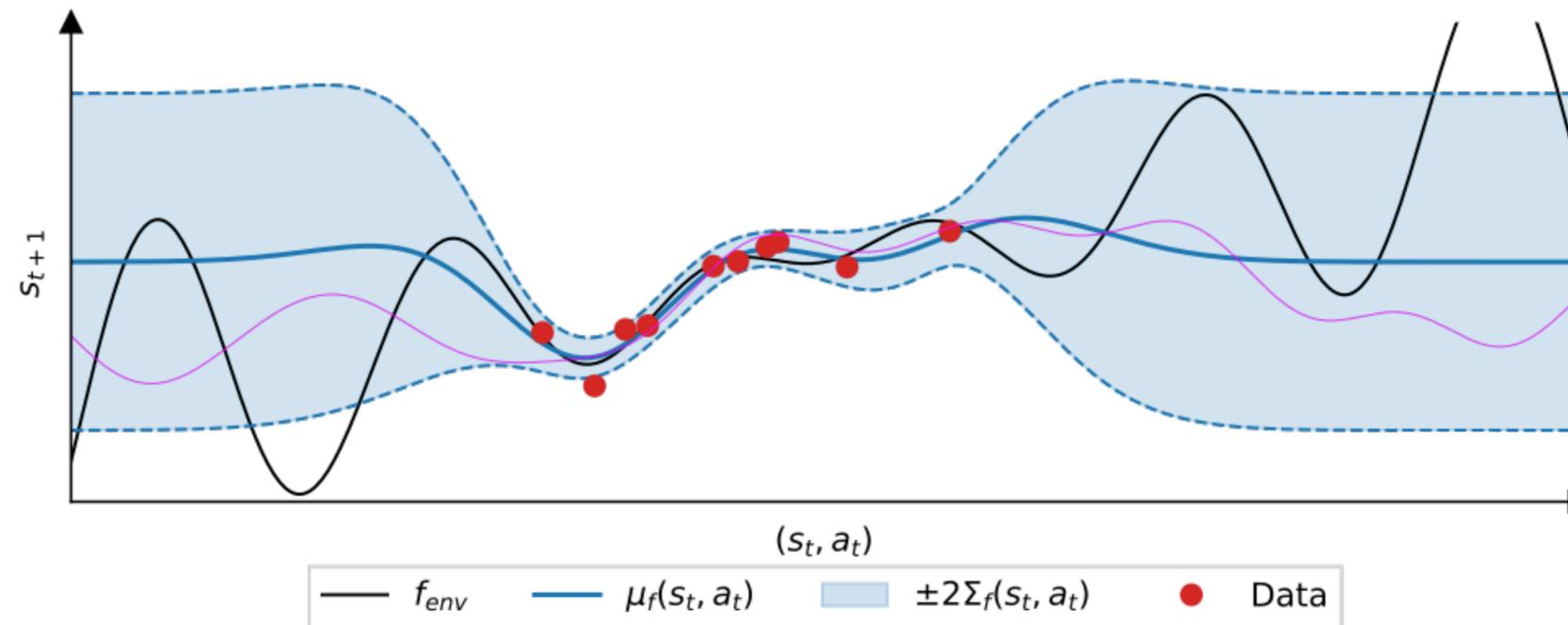
# Exploration via Posterior Sampling

$$\pi^{PS} = \arg \max_{\pi} J(\pi; \tilde{f}), \quad \tilde{f} \sim p(f | \mathcal{D})$$



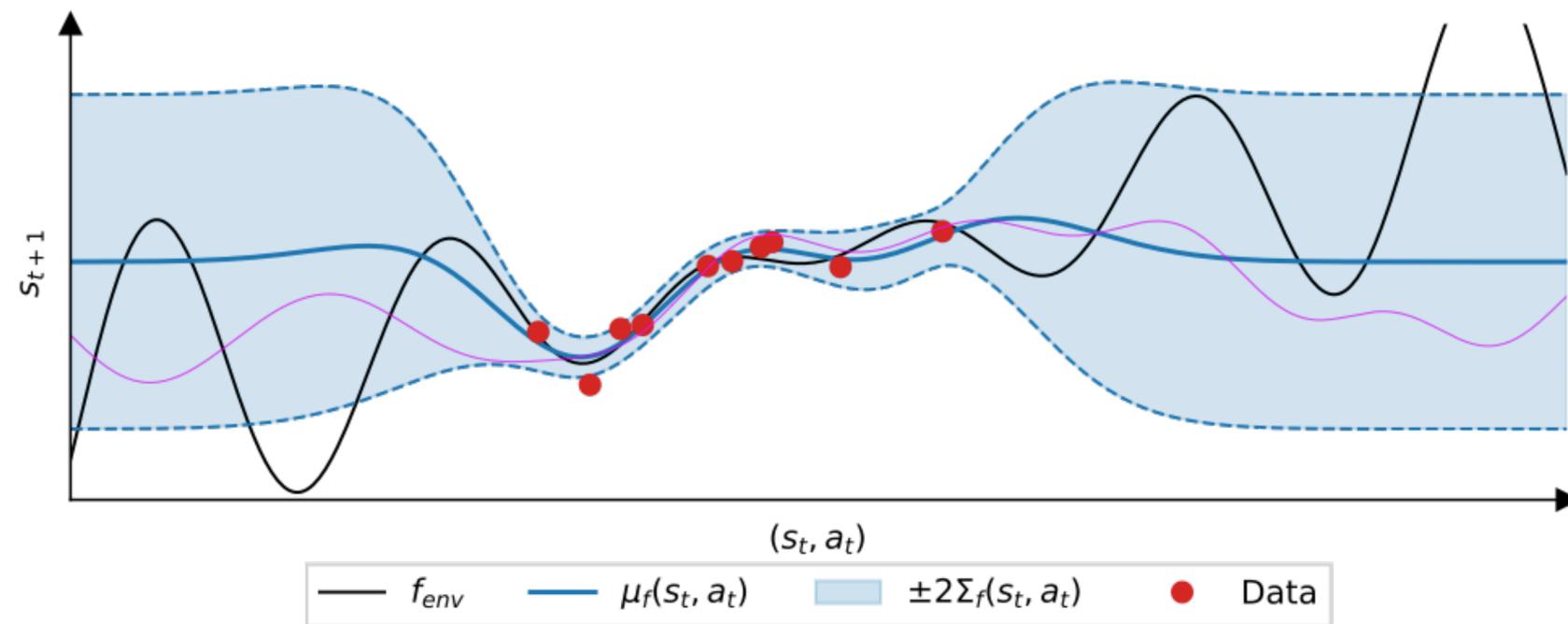
# Exploration via Posterior Sampling

$$\pi^{PS} = \arg \max_{\pi} J(\pi; \tilde{f}), \quad \tilde{f} \sim p(f | \mathcal{D})$$



# Exploration via Posterior Sampling

$$\pi^{PS} = \arg \max_{\pi} J(\pi; \tilde{f}), \quad \tilde{f} \sim p(f | \mathcal{D})$$



No extra hyperparameters

# Exploration via Upper Confidence Bound

$$\pi^{UCB} = \arg \max_{\pi} \max_{f \in \mathcal{M}} J(\pi; f)$$

# Exploration via Upper Confidence Bound

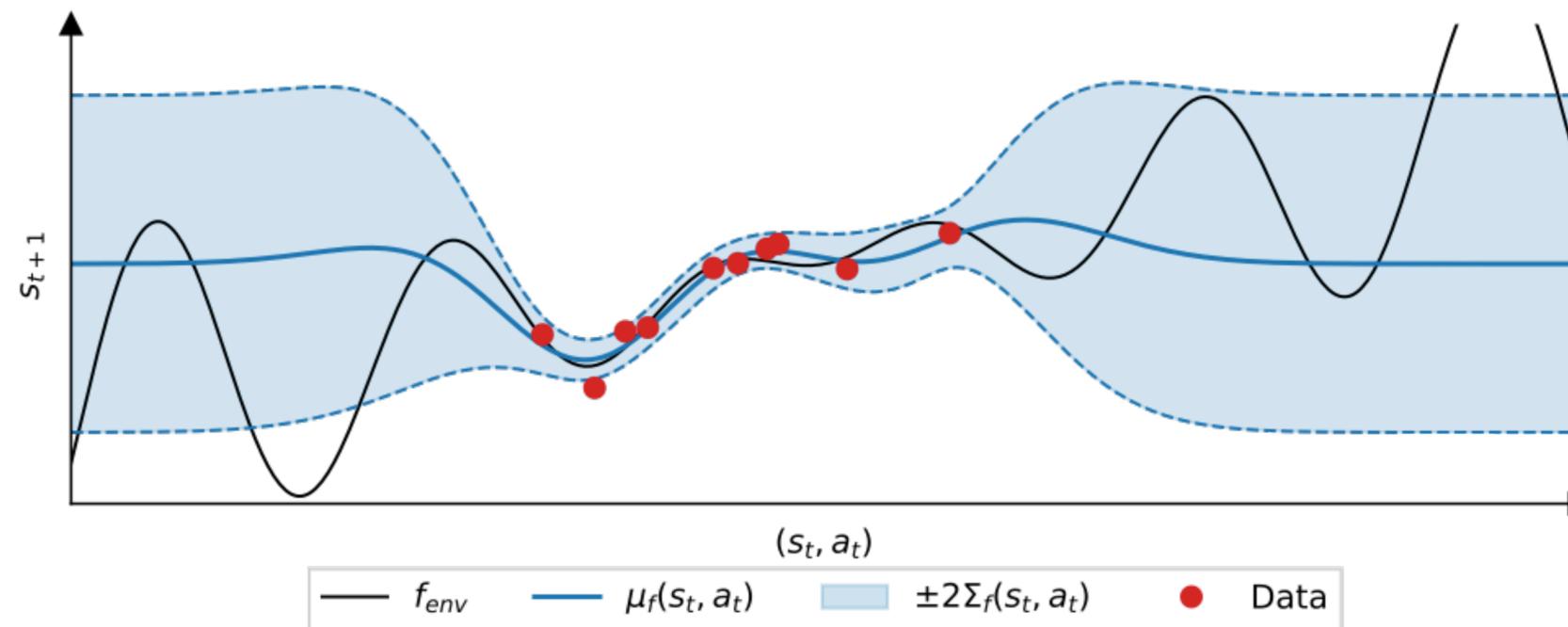
$$\pi^{UCB} = \arg \max_{\pi} \max_{f \in \mathcal{M}} J(\pi; f)$$

$$\mathcal{M} = \{f \mid \|f(s, a) - \mu_f(s, a)\| \leq \beta \Sigma_f(s, a)\}$$

# Exploration via Upper Confidence Bound

$$\pi^{UCB} = \arg \max_{\pi} \max_{f \in \mathcal{M}} J(\pi; f)$$

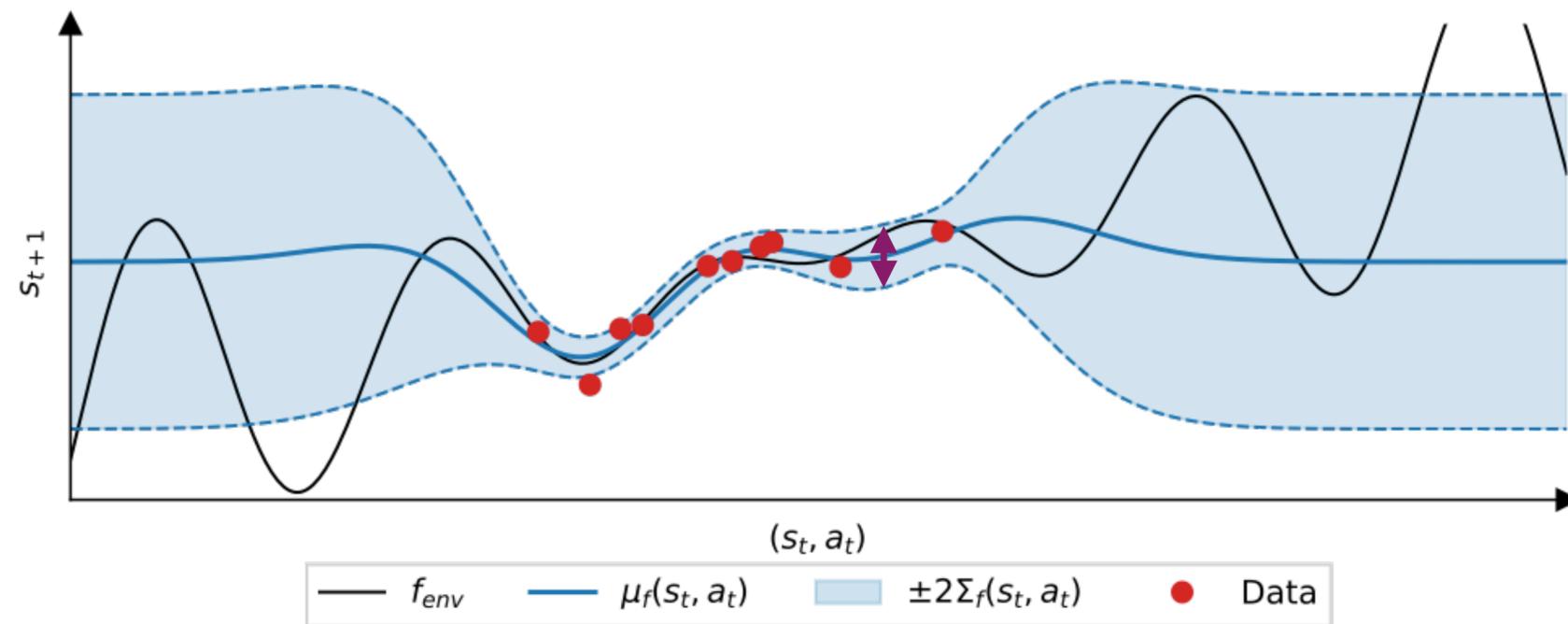
$$\mathcal{M} = \{f \mid \|f(s, a) - \mu_f(s, a)\| \leq \beta \Sigma_f(s, a)\}$$



# Exploration via Upper Confidence Bound

$$\pi^{UCB} = \arg \max_{\pi} \max_{f \in \mathcal{M}} J(\pi; f)$$

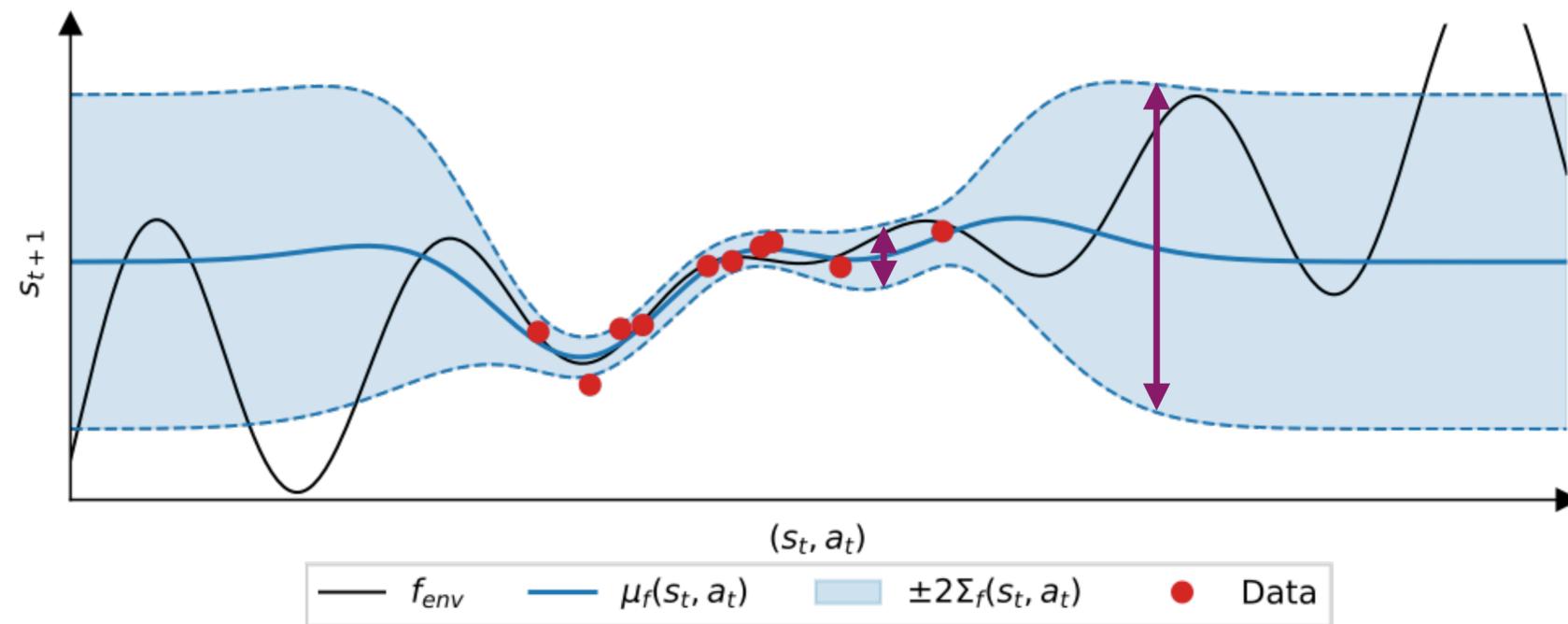
$$\mathcal{M} = \{f \mid \|f(s, a) - \mu_f(s, a)\| \leq \beta \Sigma_f(s, a)\}$$



# Exploration via Upper Confidence Bound

$$\pi^{UCB} = \arg \max_{\pi} \max_{f \in \mathcal{M}} J(\pi; f)$$

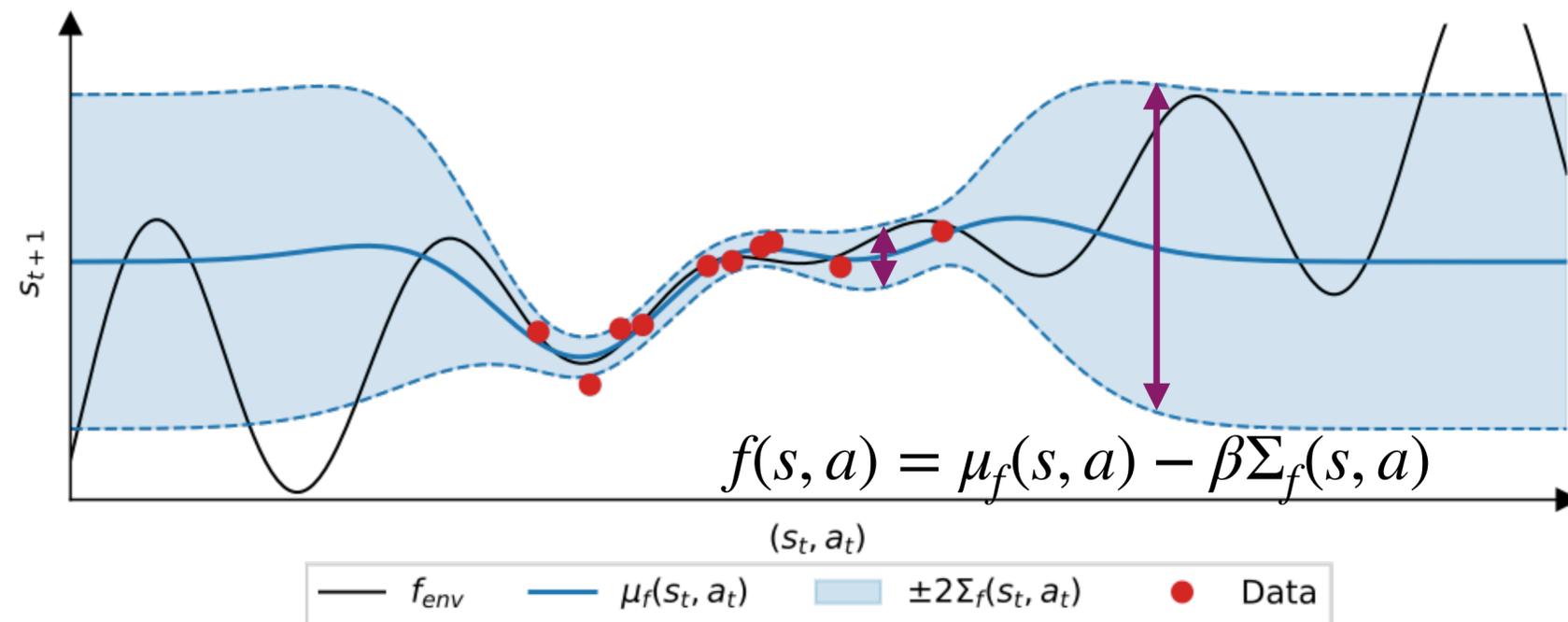
$$\mathcal{M} = \{f \mid \|f(s, a) - \mu_f(s, a)\| \leq \beta \Sigma_f(s, a)\}$$



# Exploration via Upper Confidence Bound

$$\pi^{UCB} = \arg \max_{\pi} \max_{f \in \mathcal{M}} J(\pi; f)$$

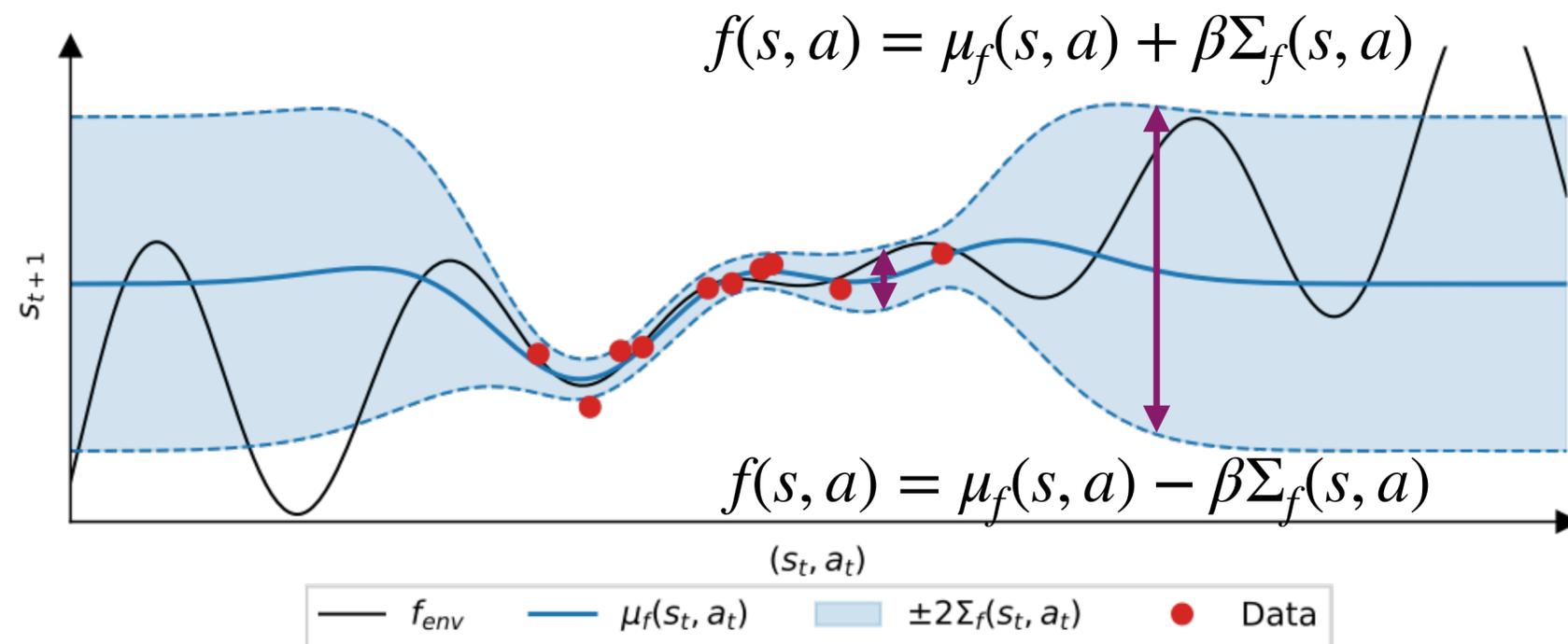
$$\mathcal{M} = \{f \mid \|f(s, a) - \mu_f(s, a)\| \leq \beta \Sigma_f(s, a)\}$$



# Exploration via Upper Confidence Bound

$$\pi^{UCB} = \arg \max_{\pi} \max_{f \in \mathcal{M}} J(\pi; f)$$

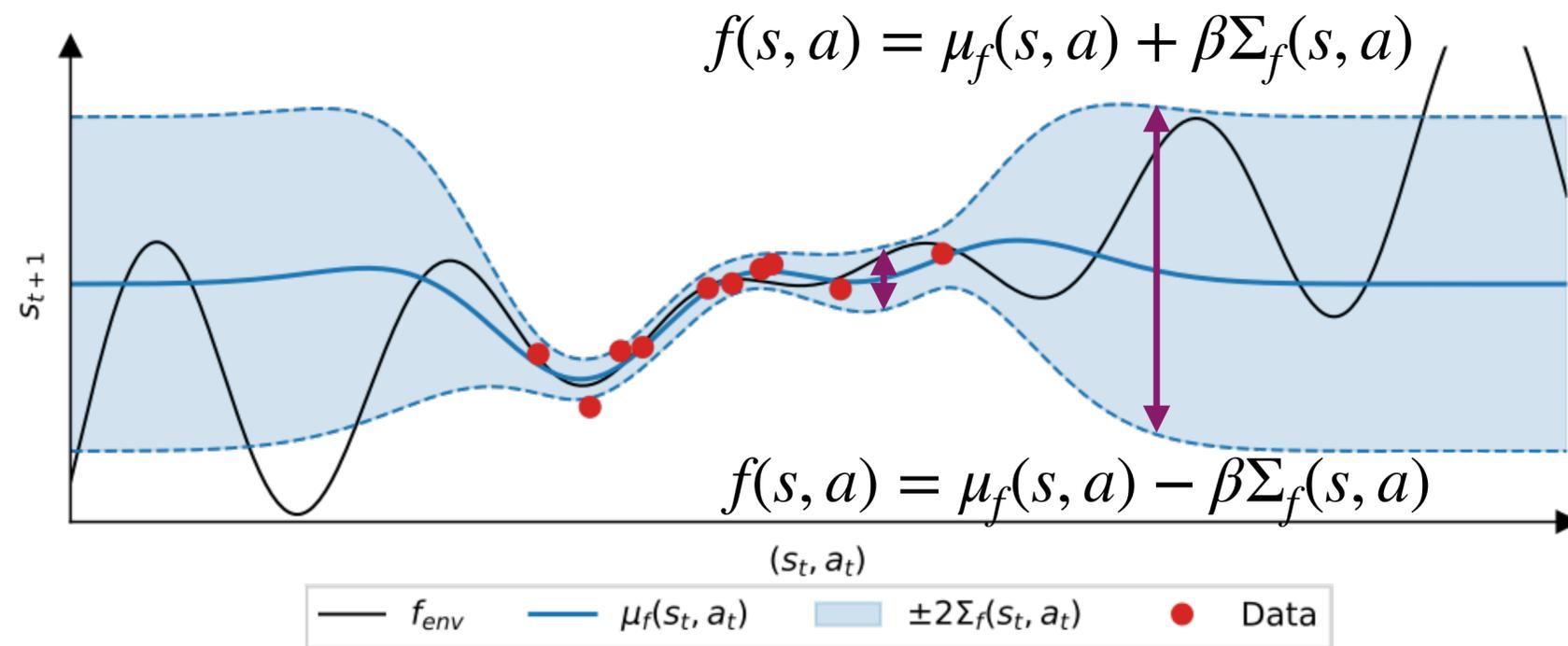
$$\mathcal{M} = \{f \mid \|f(s, a) - \mu_f(s, a)\| \leq \beta \Sigma_f(s, a)\}$$



# Exploration via Upper Confidence Bound

$$\pi^{UCB} = \arg \max_{\pi} \max_{f \in \mathcal{M}} J(\pi; f)$$

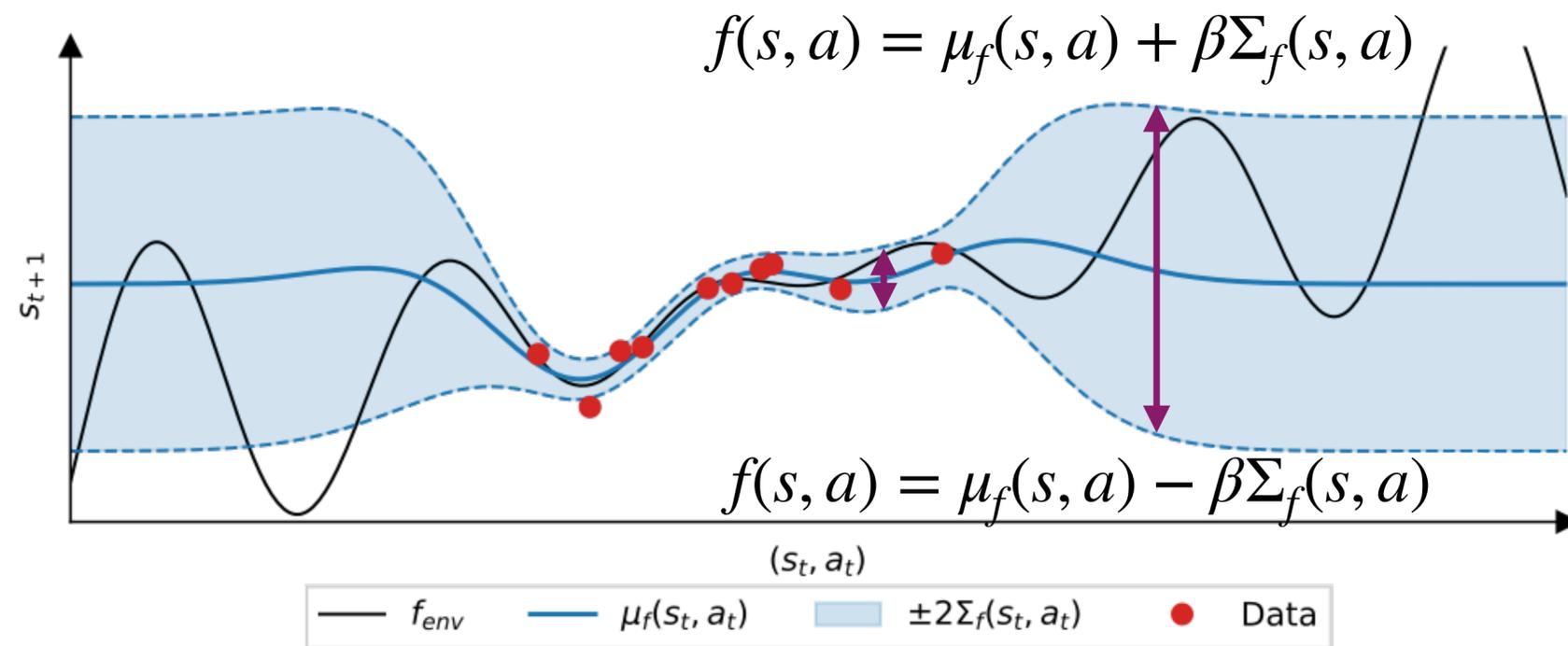
$$\mathcal{M} = \{f \mid \|f(s, a) - \mu_f(s, a)\| \leq \beta \Sigma_f(s, a)\}$$



# Exploration via Upper Confidence Bound

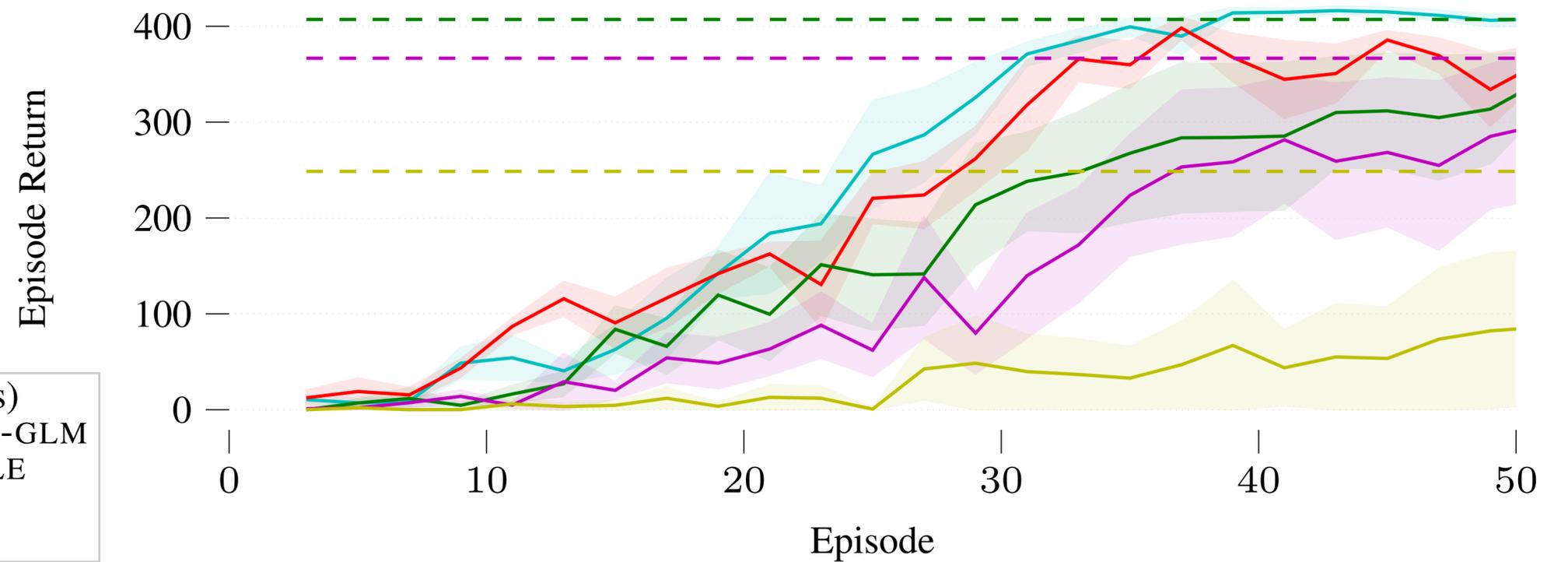
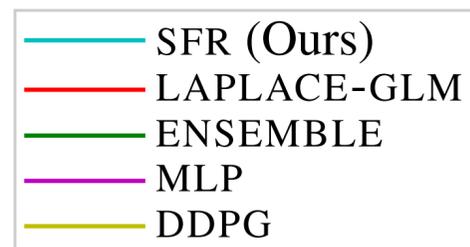
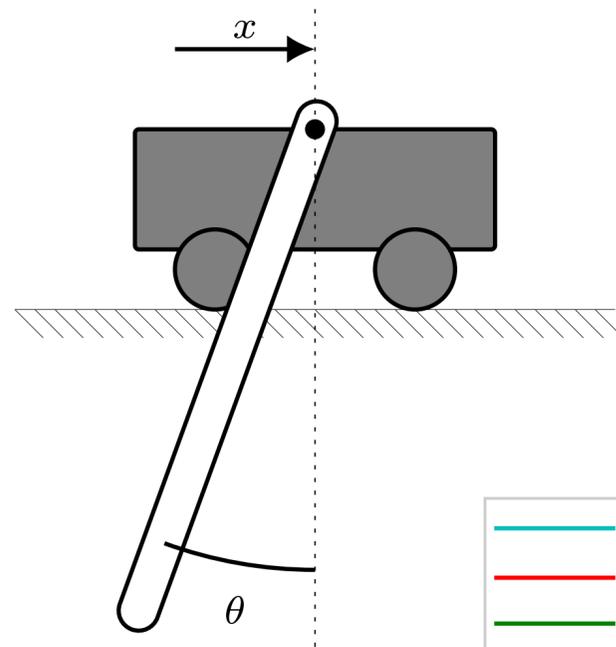
$$\pi^{UCB} = \arg \max_{\pi} \max_{f \in \mathcal{M}} J(\pi; f)$$

$$\mathcal{M} = \{f \mid \|f(s, a) - \mu_f(s, a)\| \leq \beta \Sigma_f(s, a)\}$$



Extra hyperparameter  $\beta$

# How to Quantify Uncertainty in Dynamics?



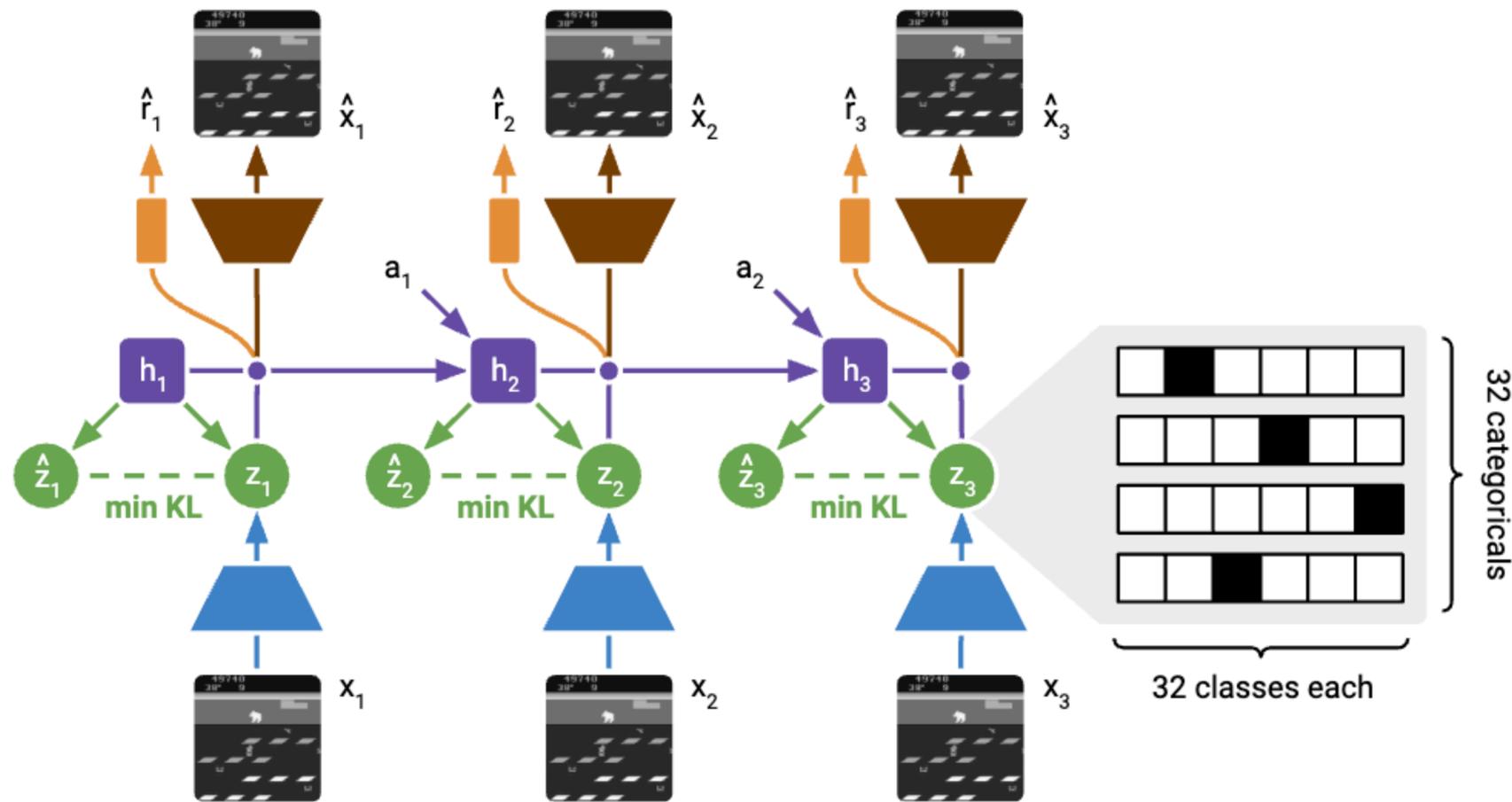
# Learning Objectives

Understand

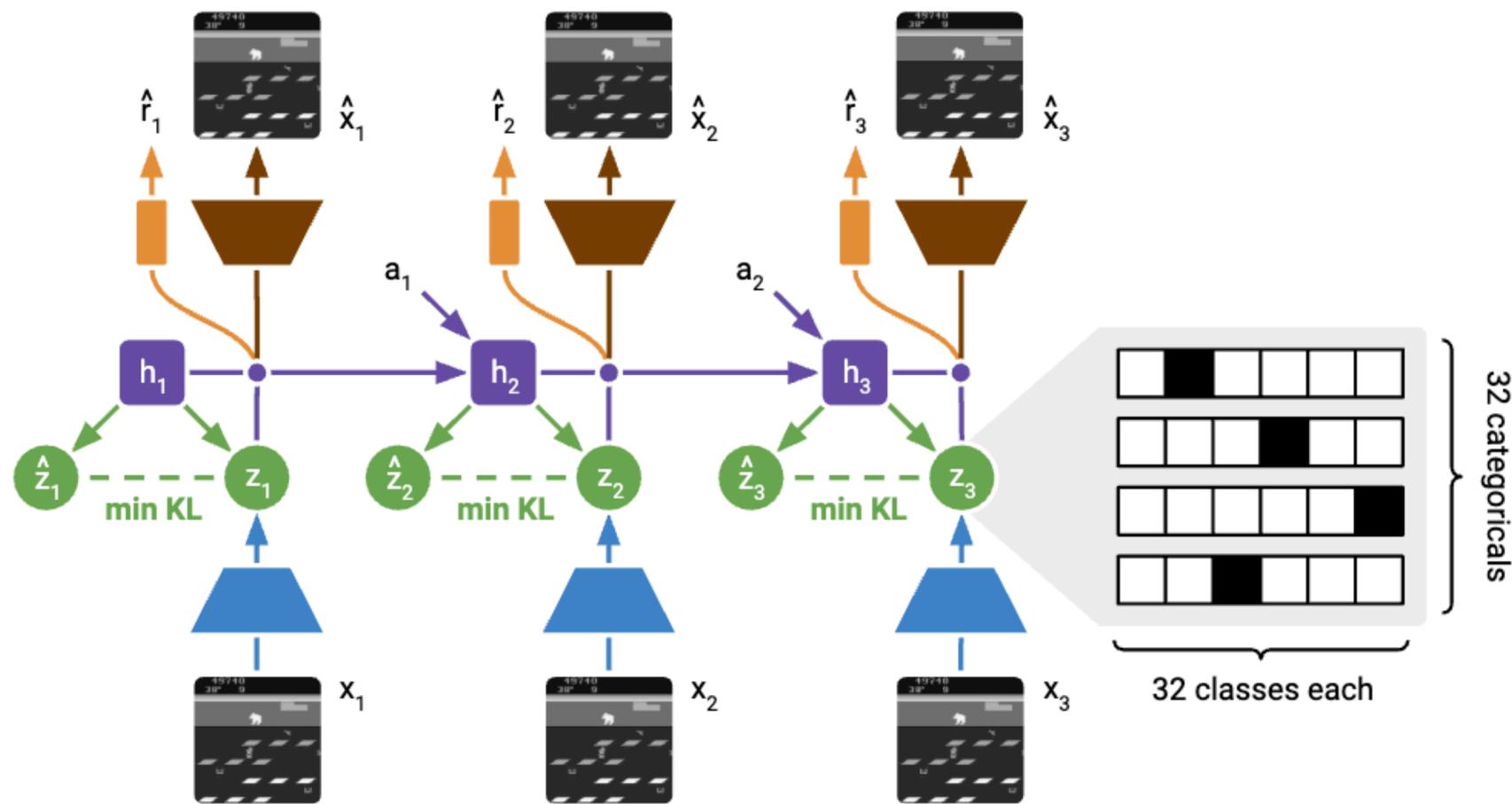
1. ~~What a “model” is in model-based RL~~
2. ~~How a “model” can aid decision making~~
3. ~~Differences between background and decision-time planning~~
4. ~~Decision-time planning strategies for continuous actions~~
5. ~~Sources of uncertainty in model-based RL~~
6. ~~Rationale and insights for decision-making under uncertainty~~

# Case Study: Dreamer (v2/v3)

World Model Learning



# Case Study: Dreamer (v2/v3)

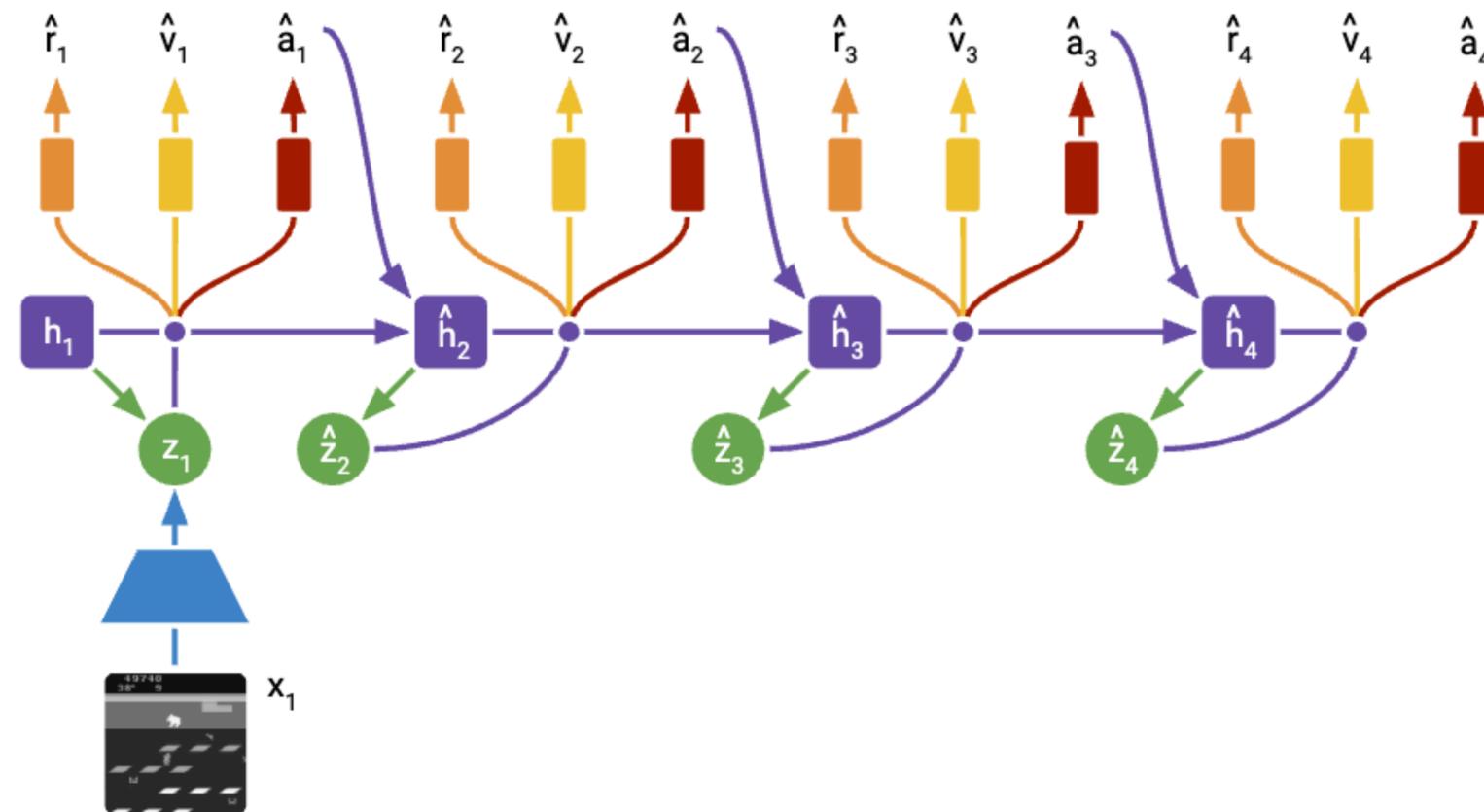


World Model Learning

Latent dynamics with encoder/decoder

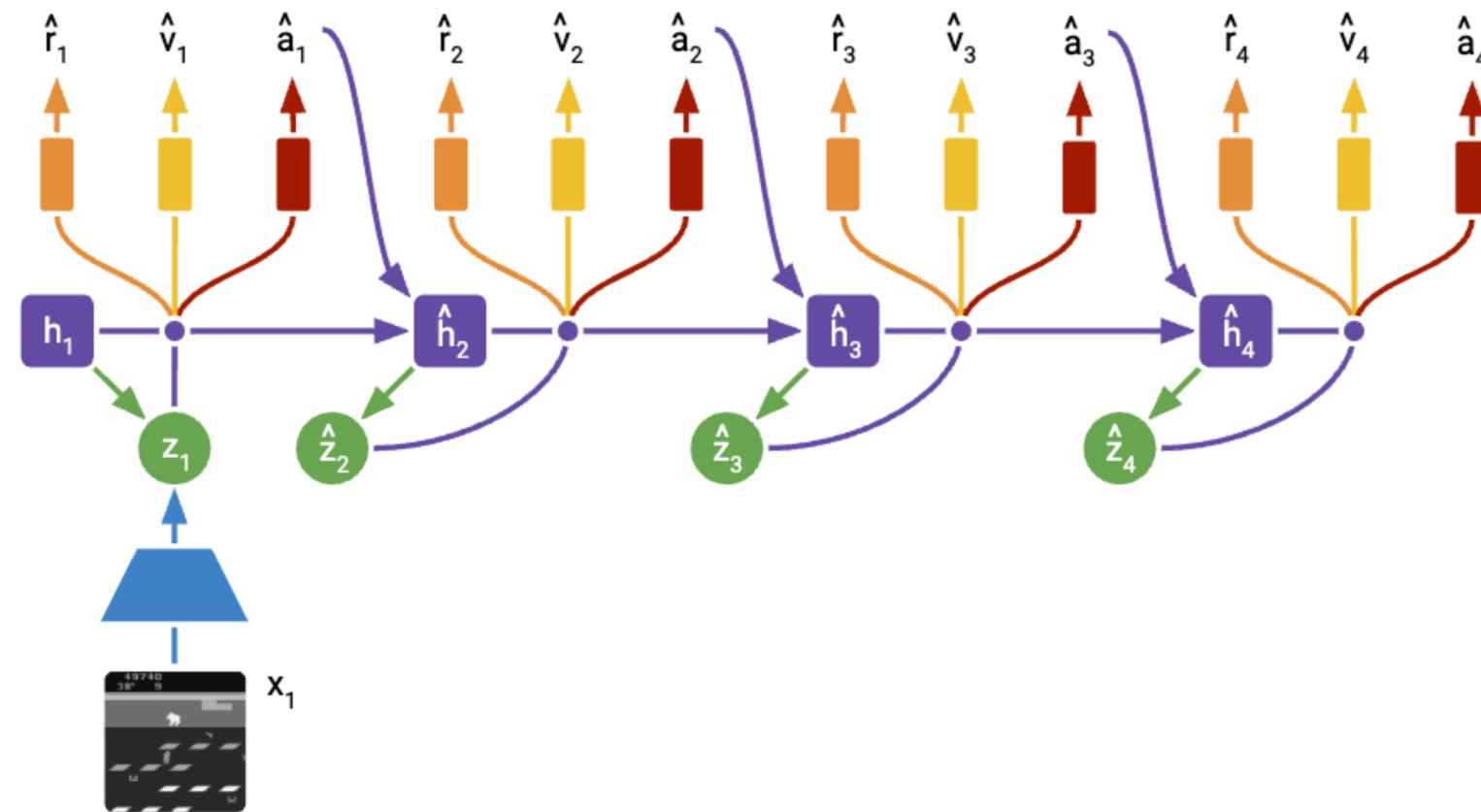
# Case Study: Dreamer (v2/v3)

Actor Critic Learning



Latent dynamics with encoder/decoder

# Case Study: Dreamer (v2/v3)

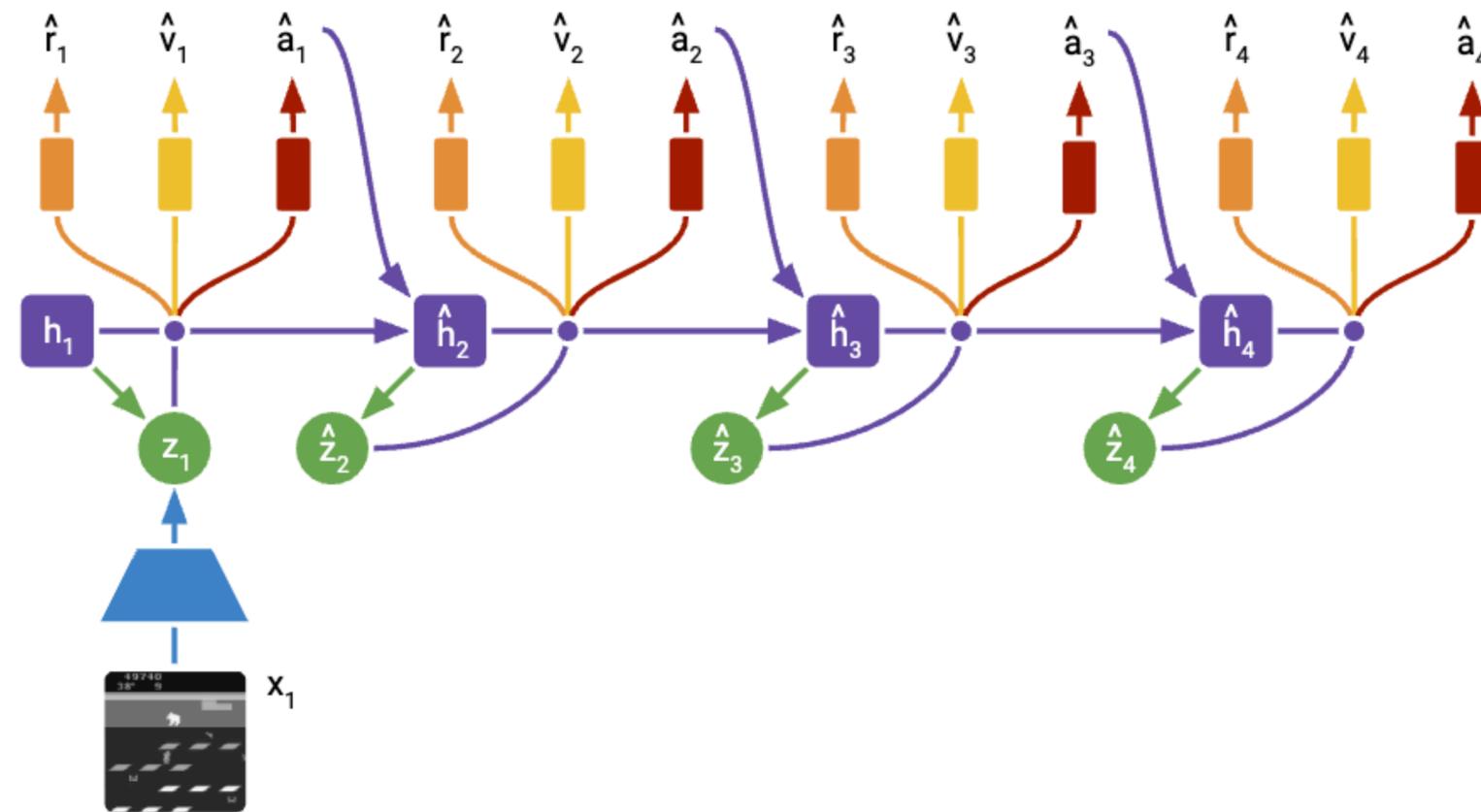


## Actor Critic Learning

Latent dynamics with encoder/decoder

Actor  $\pi_\theta(z)$  & critic  $Q_\theta(z, a)$  in latent space

# Case Study: Dreamer (v2/v3)



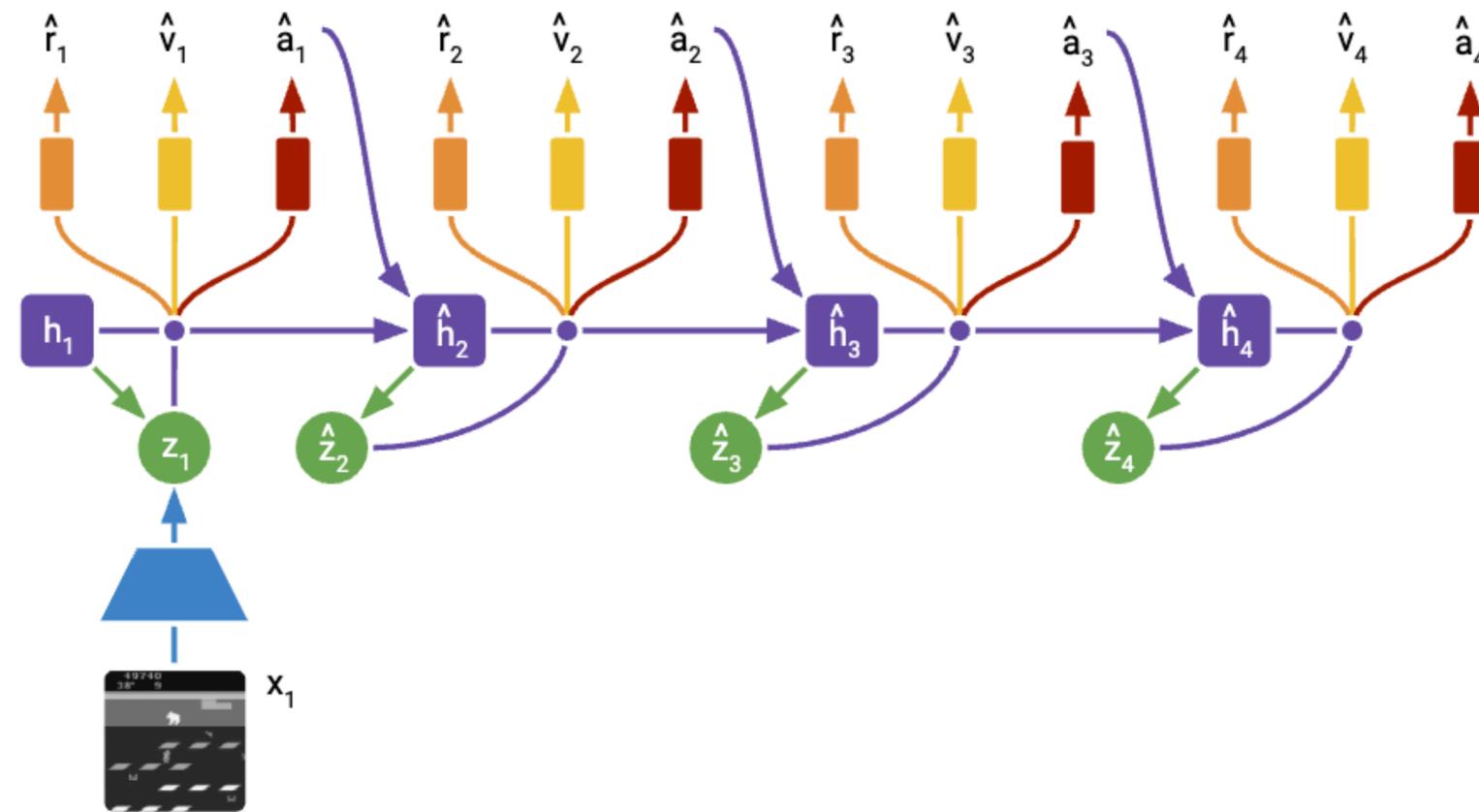
## Actor Critic Learning

Latent dynamics with encoder/decoder

Actor  $\pi_\theta(z)$  & critic  $Q_\theta(z, a)$  in latent space

Actor/critic leverage “imagined” outcomes

# Case Study: Dreamer (v2/v3)



## Actor Critic Learning

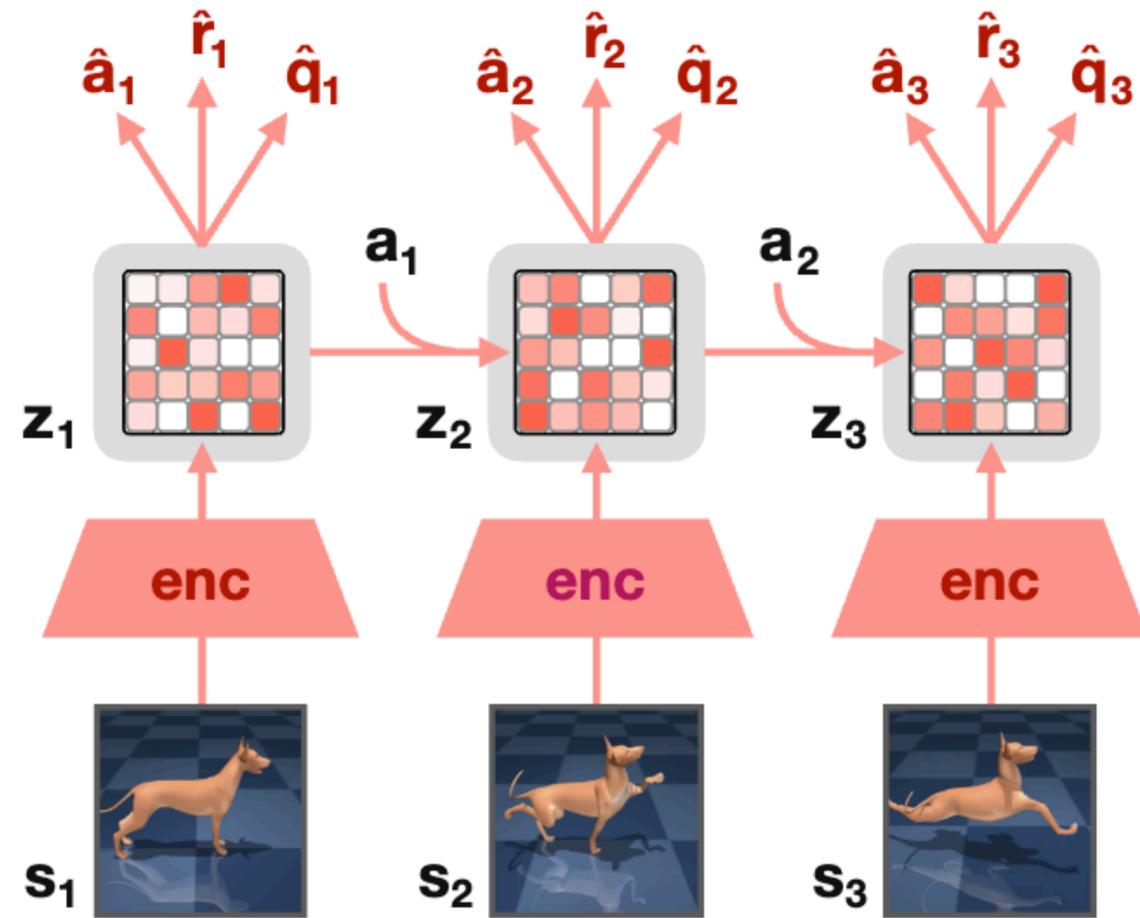
Latent dynamics with encoder/decoder

Actor  $\pi_\theta(z)$  & critic  $Q_\theta(z, a)$  in latent space

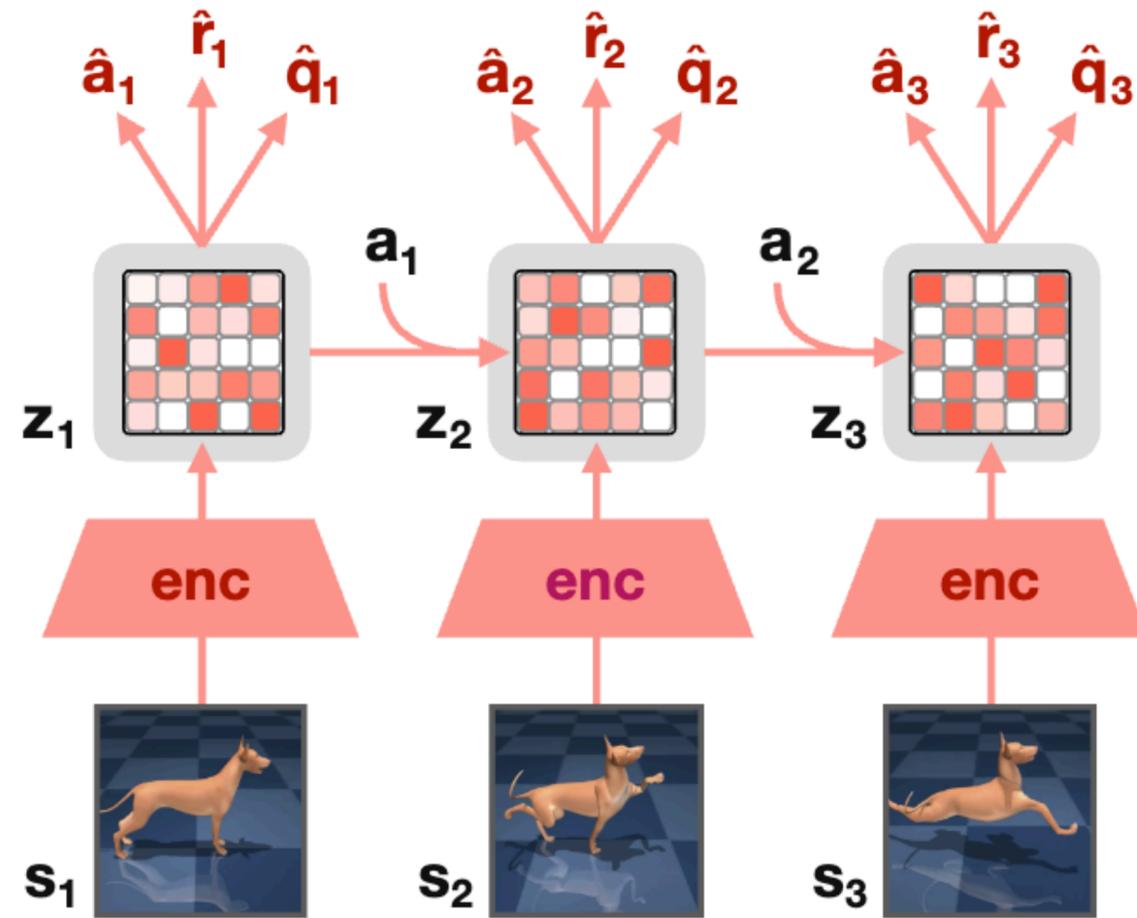
Actor/critic leverage “imagined” outcomes

Background planning

# Case Study: TD-MPC2

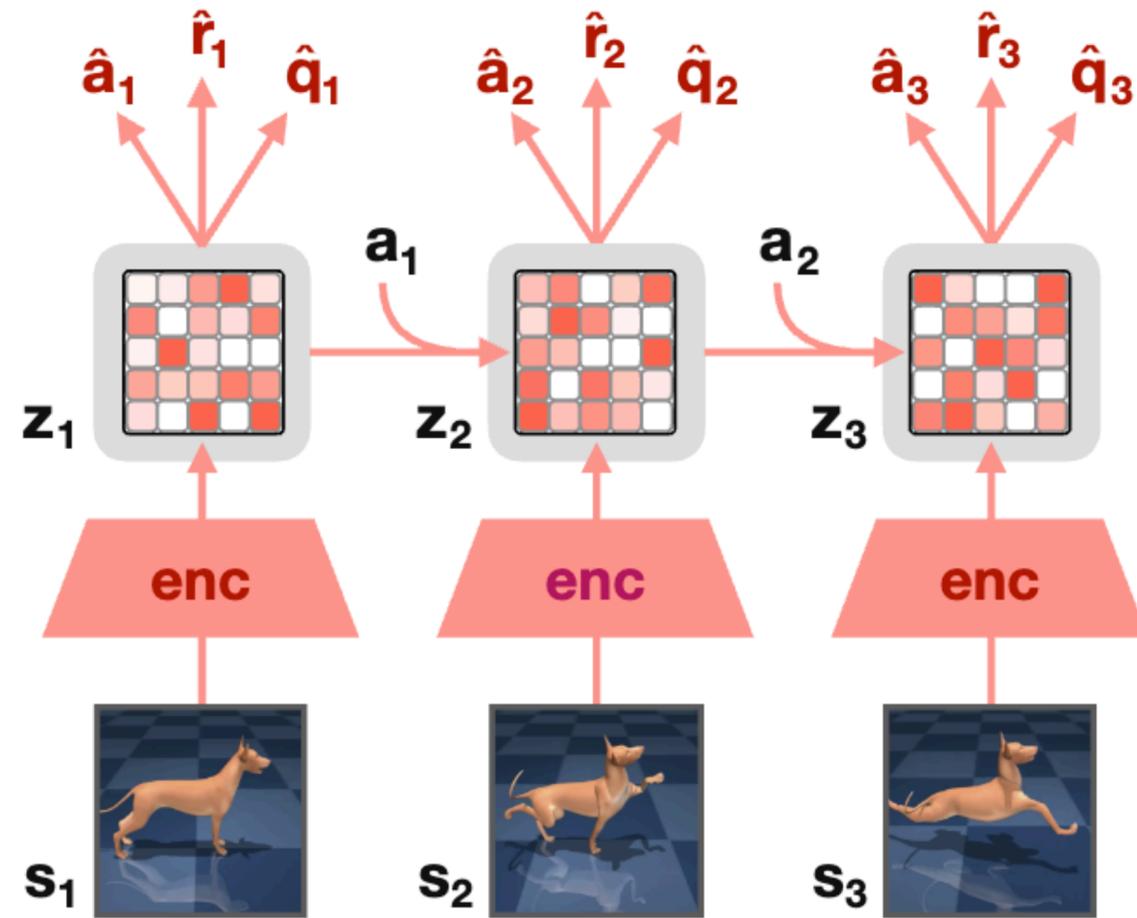


# Case Study: TD-MPC2



Latent dynamics

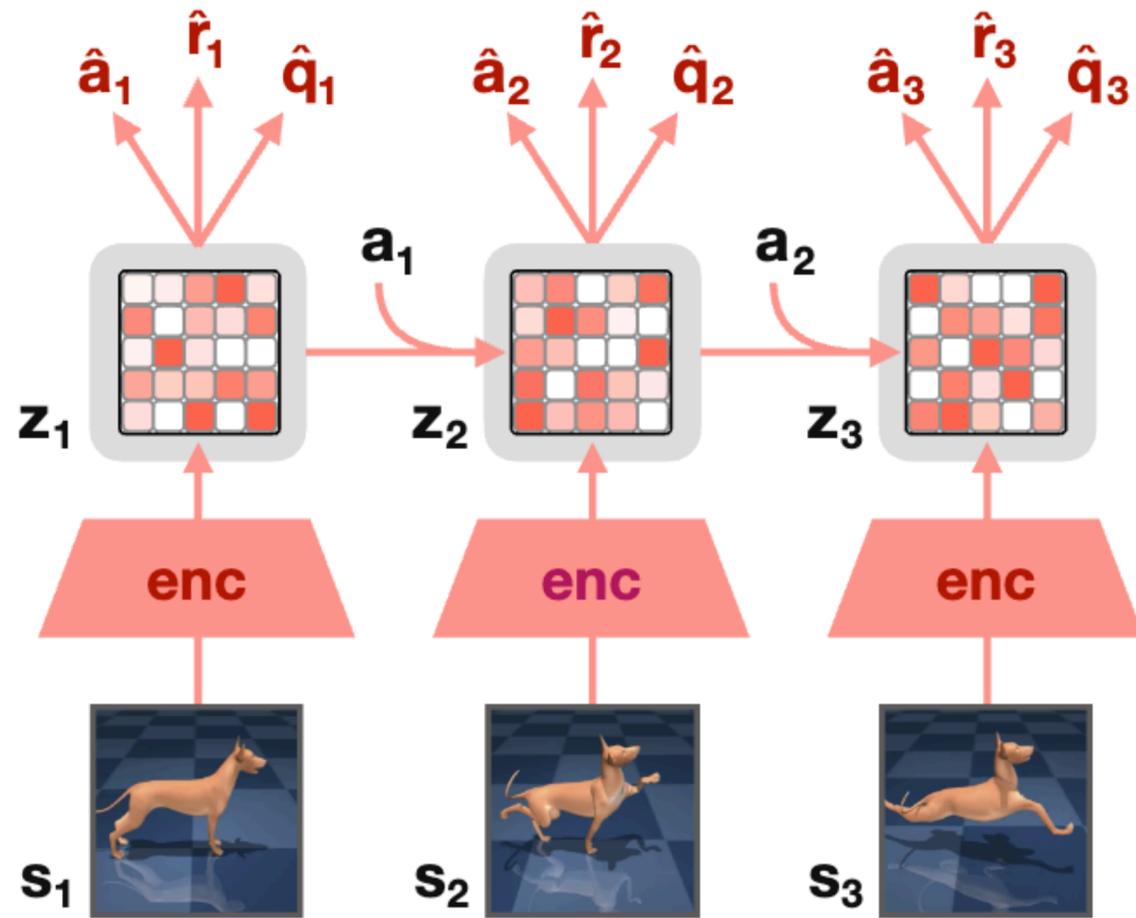
# Case Study: TD-MPC2



Latent dynamics

No decoder

# Case Study: TD-MPC2

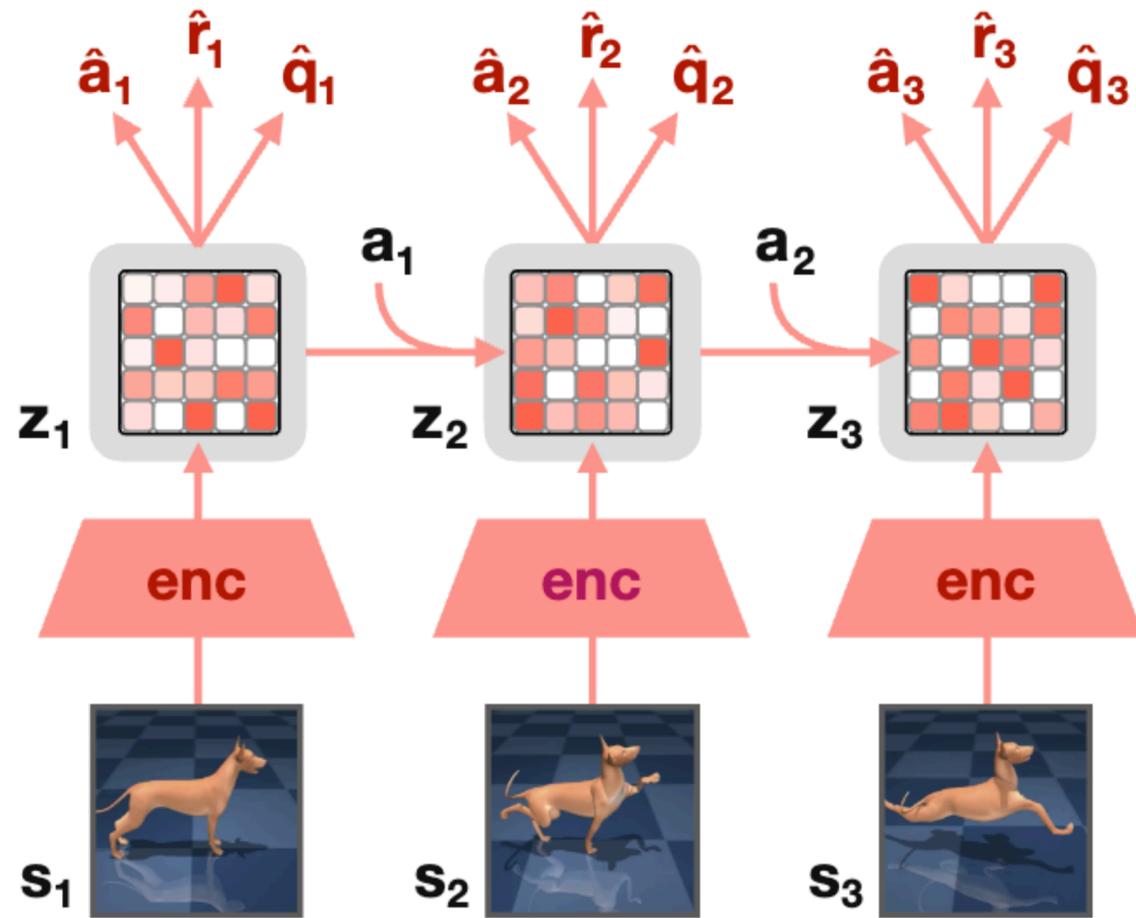


Latent dynamics

No decoder

Actor  $\pi_\theta(z)$  & critic  $Q_\theta(z, a)$  in latent space

# Case Study: TD-MPC2



Latent dynamics

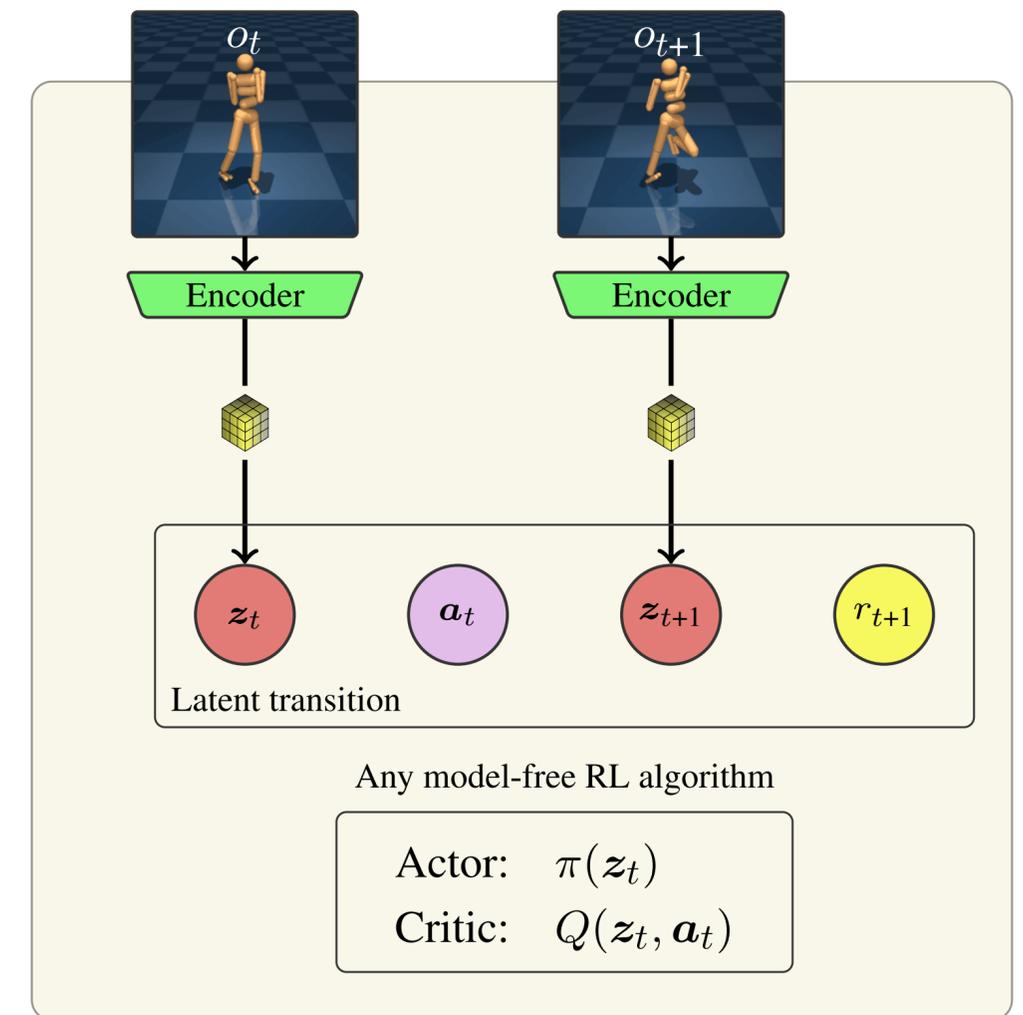
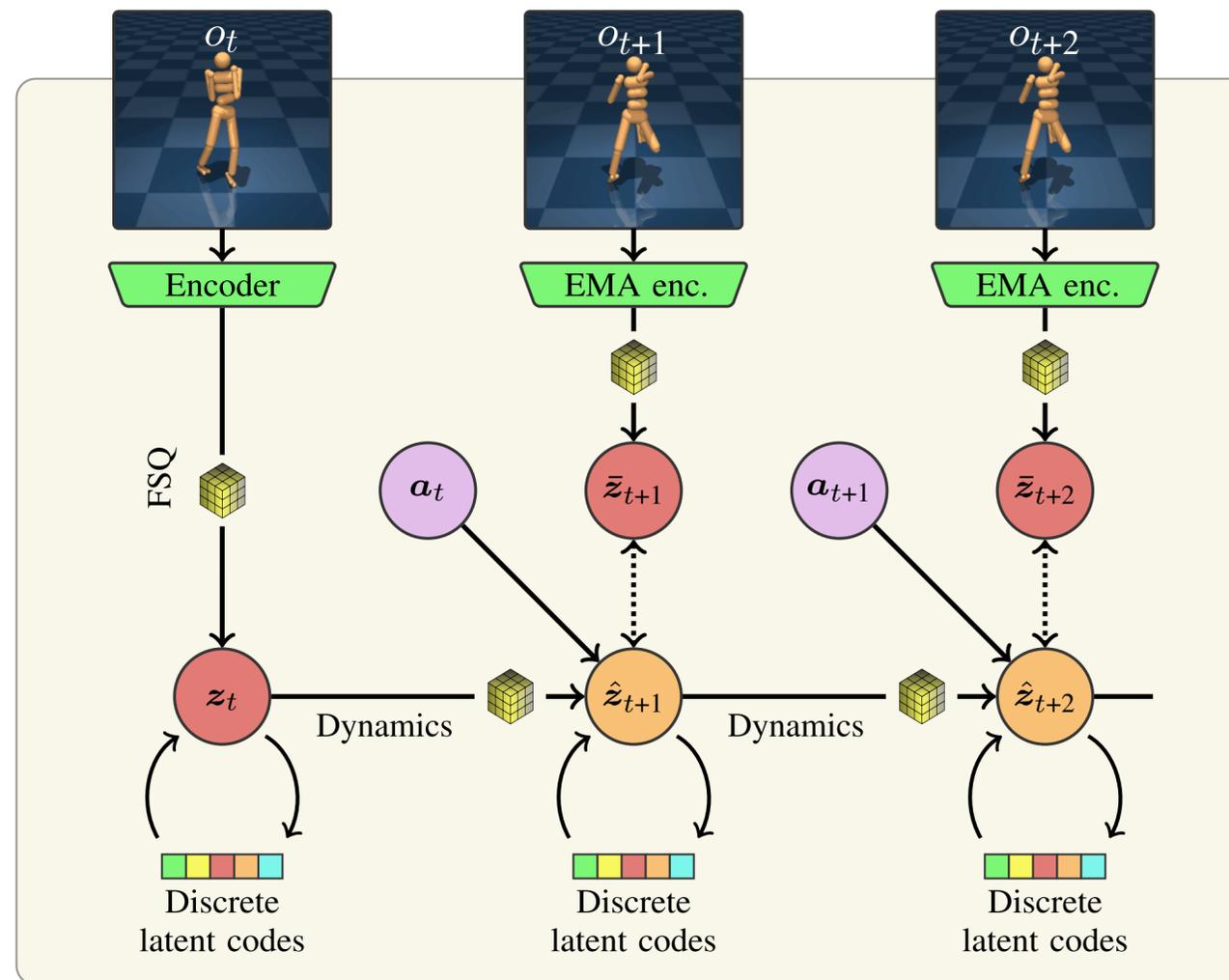
No decoder

Actor  $\pi_\theta(z)$  & critic  $Q_\theta(z, a)$  in latent space

Decision-time planning

# Case Study: iQRL

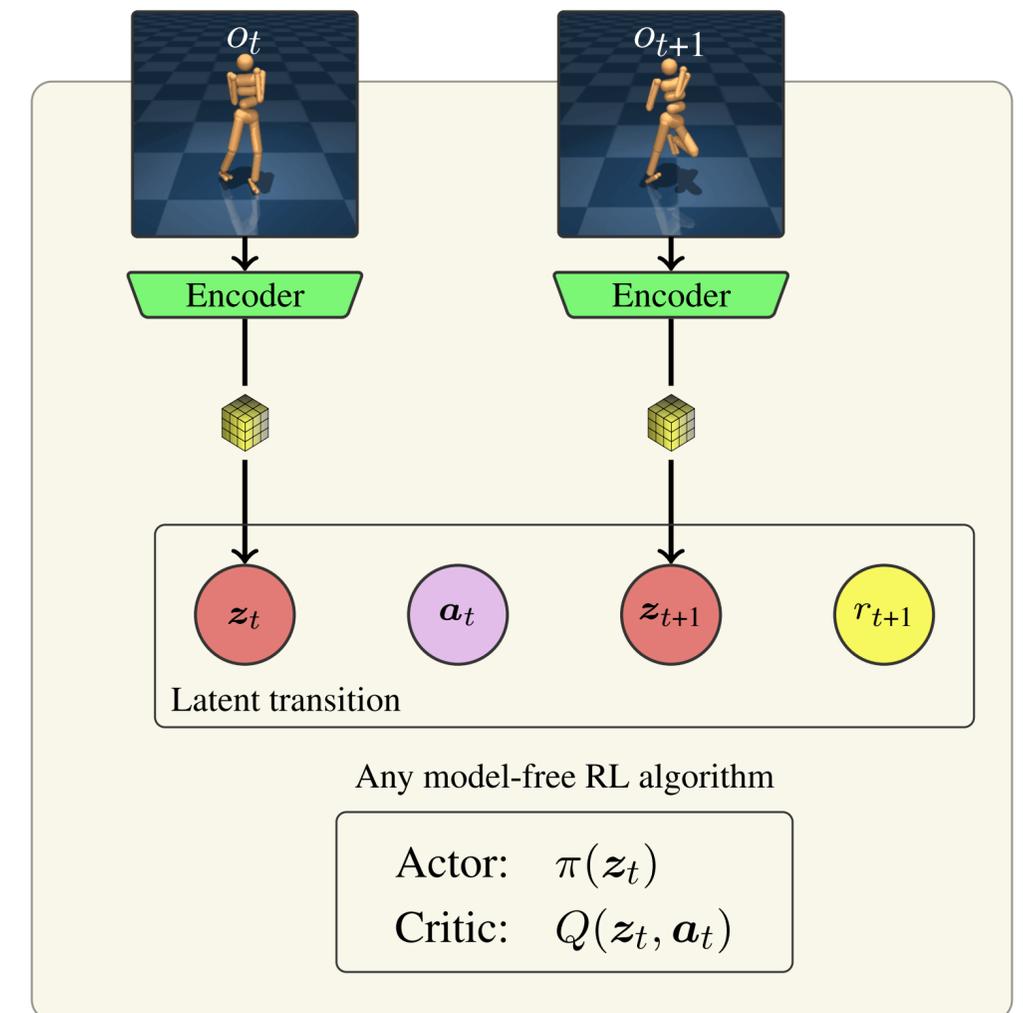
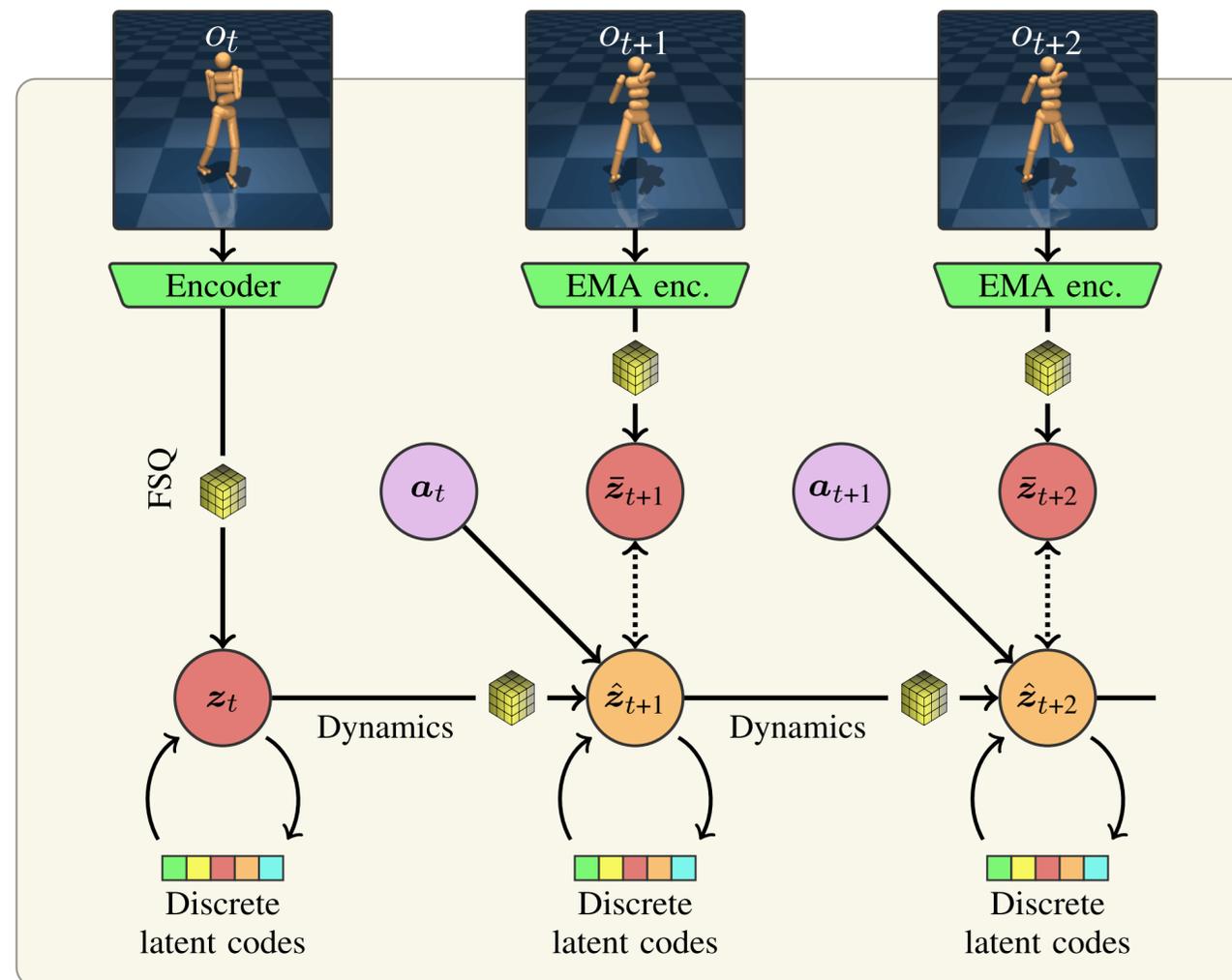
Dynamics Model but Not Model-based RL? 🤔



# Case Study: iQRL

Dynamics Model but Not Model-based RL? 🤔

Use dynamics for representation learning

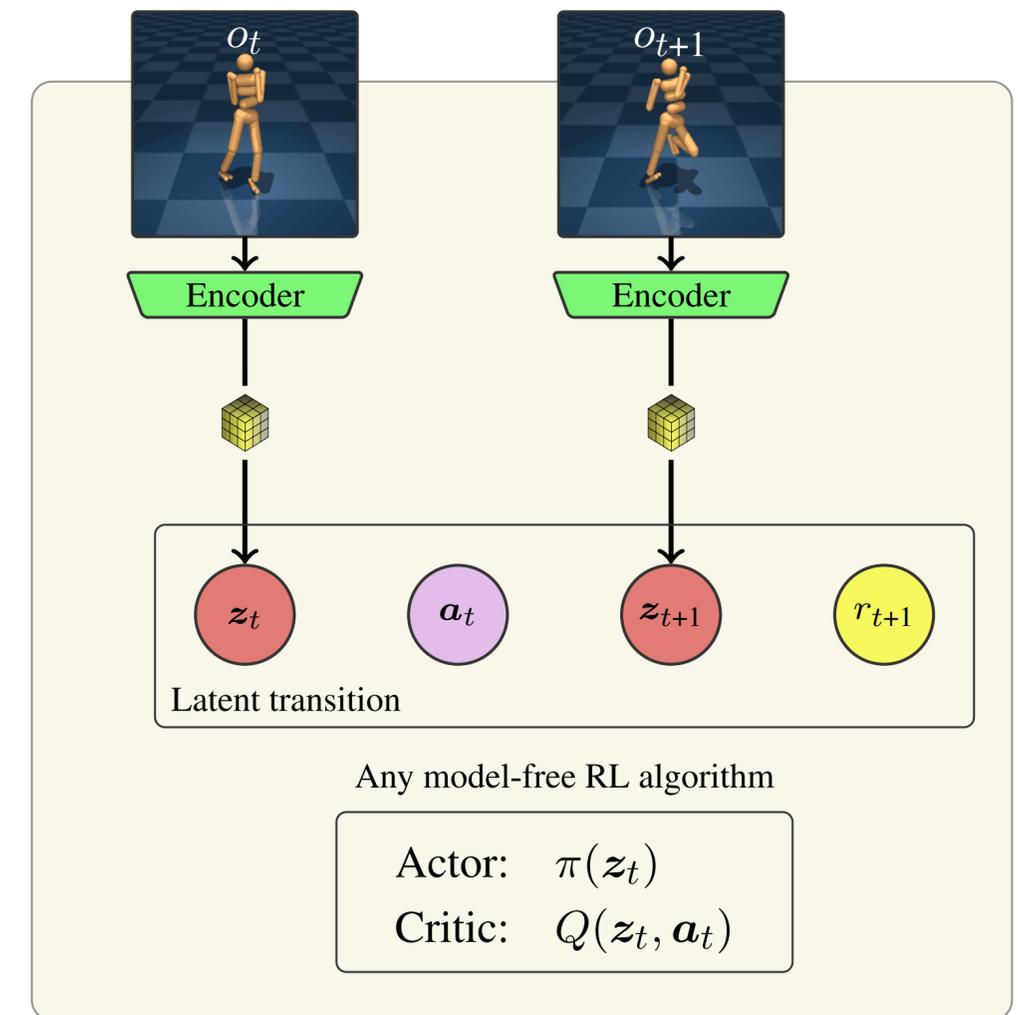
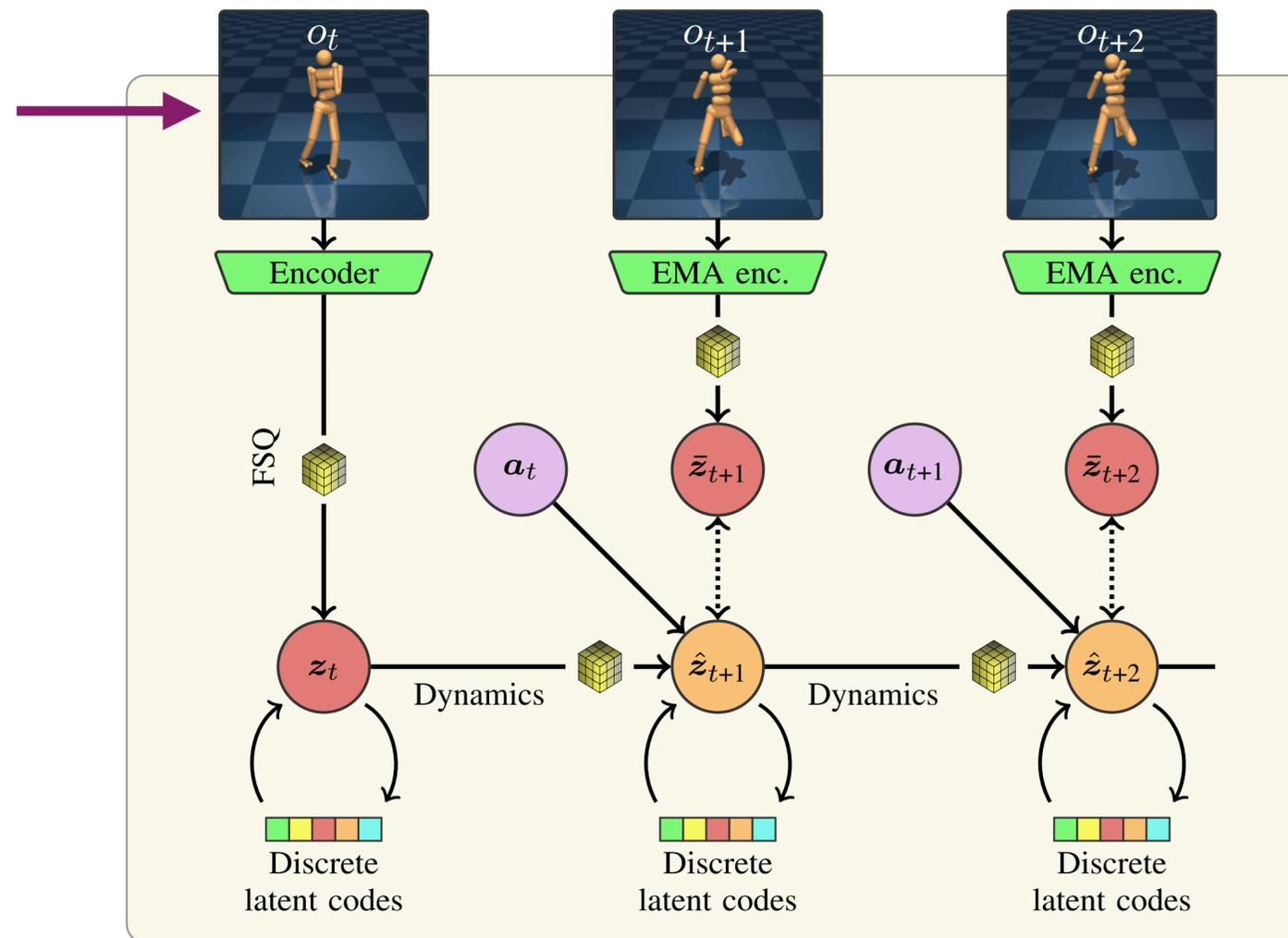


# Case Study: iQRL

Dynamics Model but Not Model-based RL? 🤔

Use dynamics for representation learning

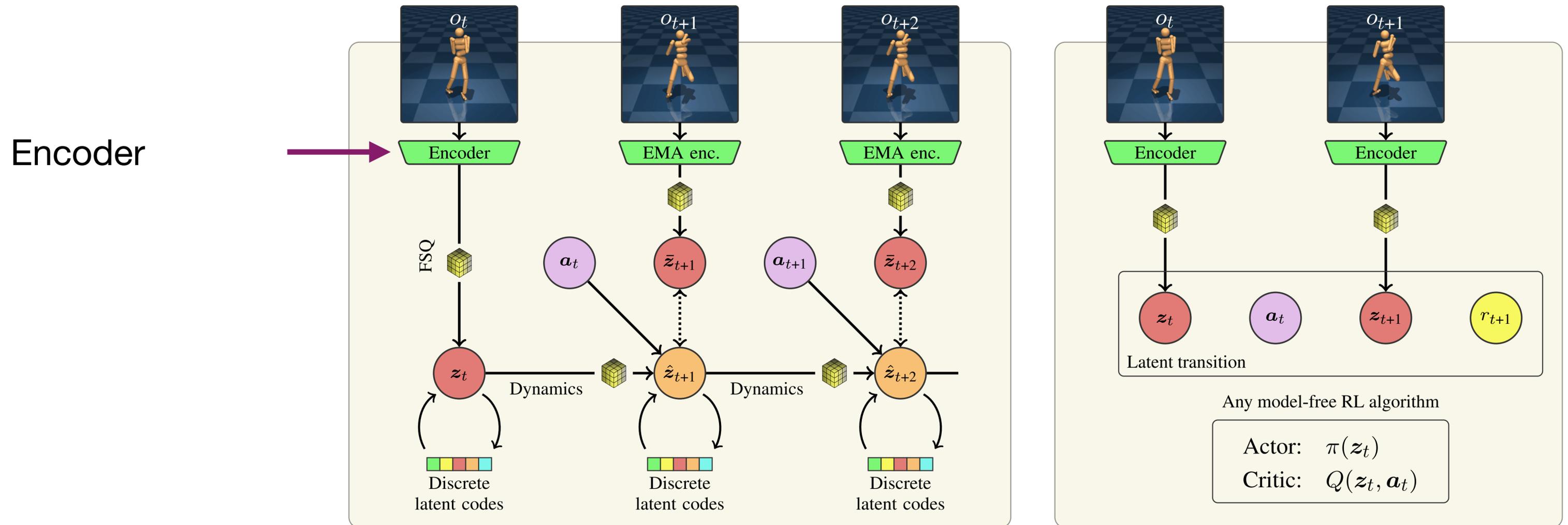
Observation



# Case Study: iQRL

Dynamics Model but Not Model-based RL? 🤔

Use dynamics for representation learning

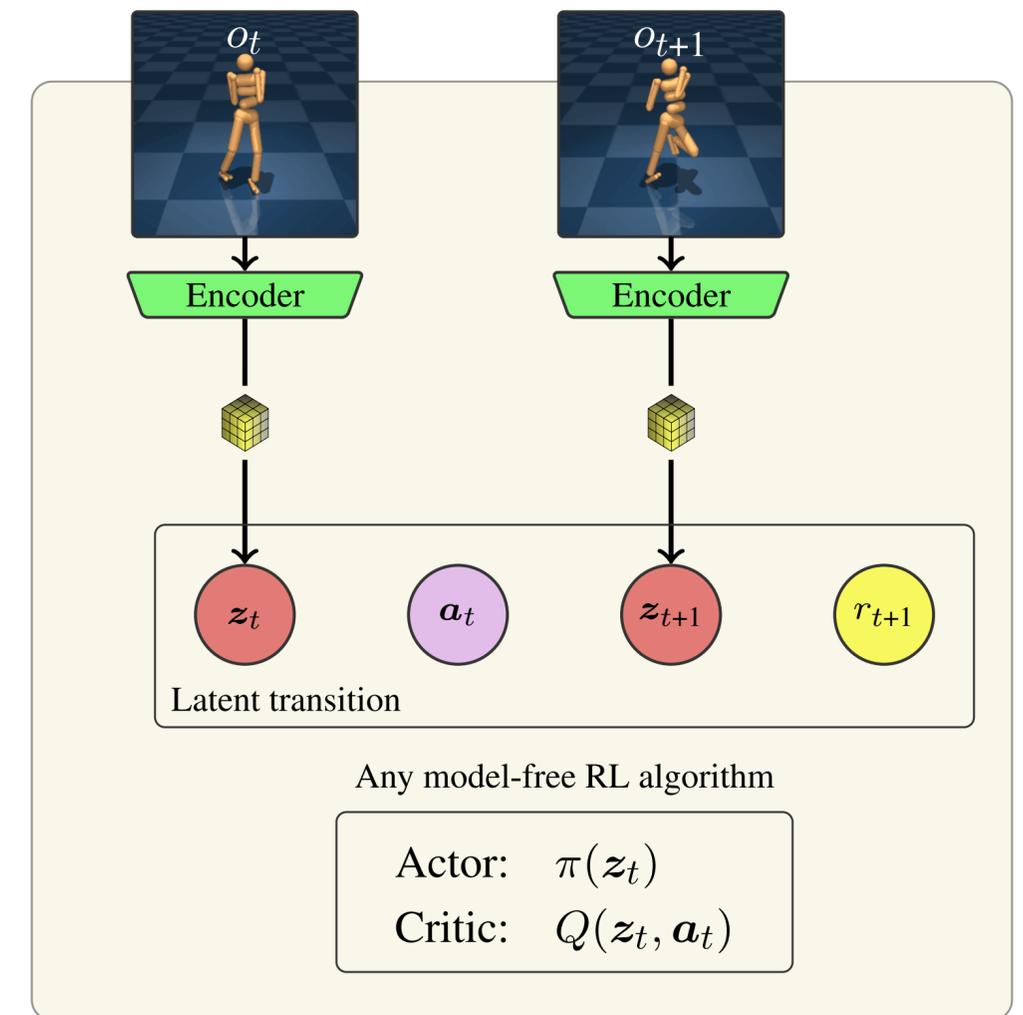
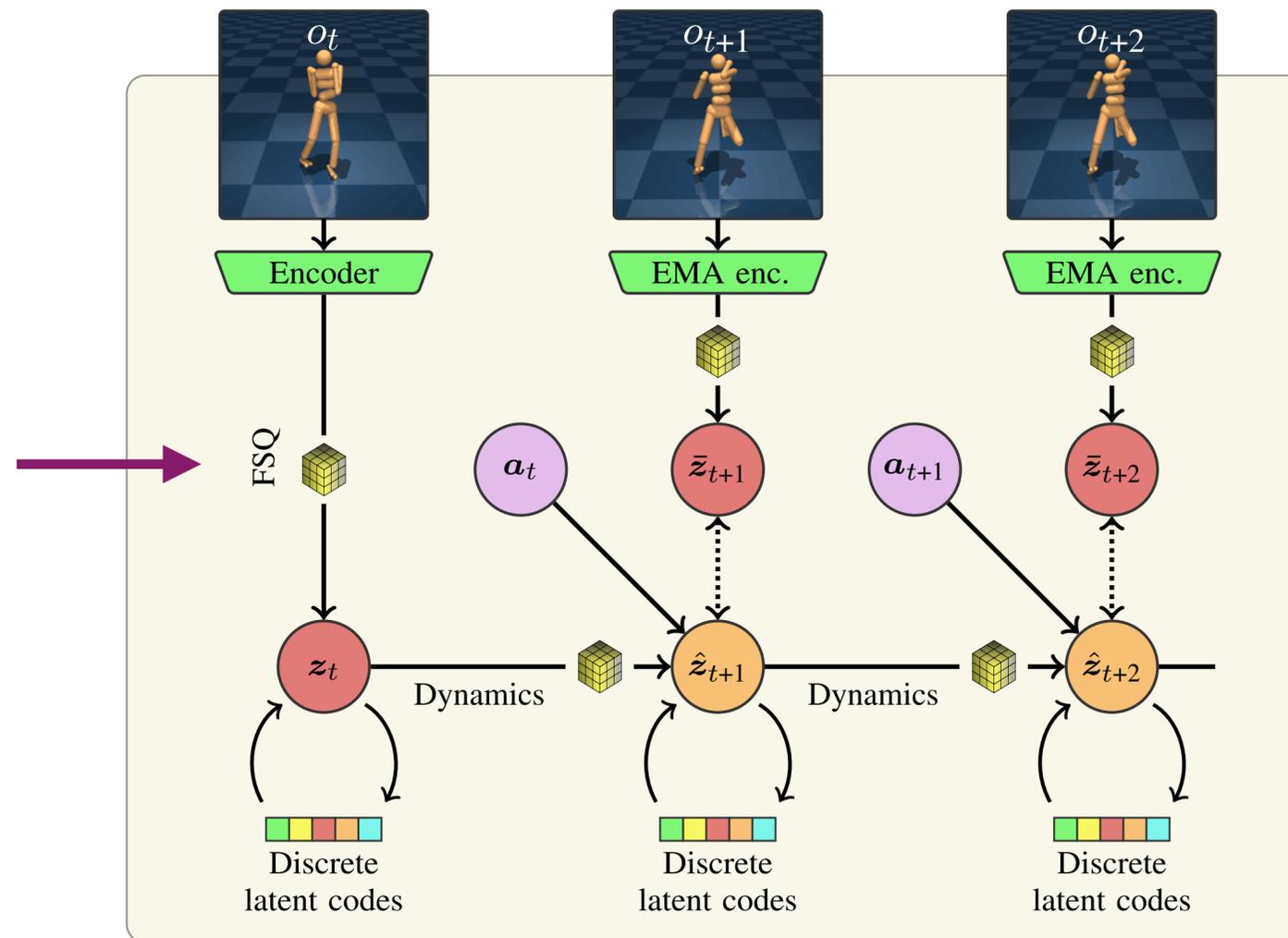


# Case Study: iQRL

Dynamics Model but Not Model-based RL? 🤔

Use dynamics for representation learning

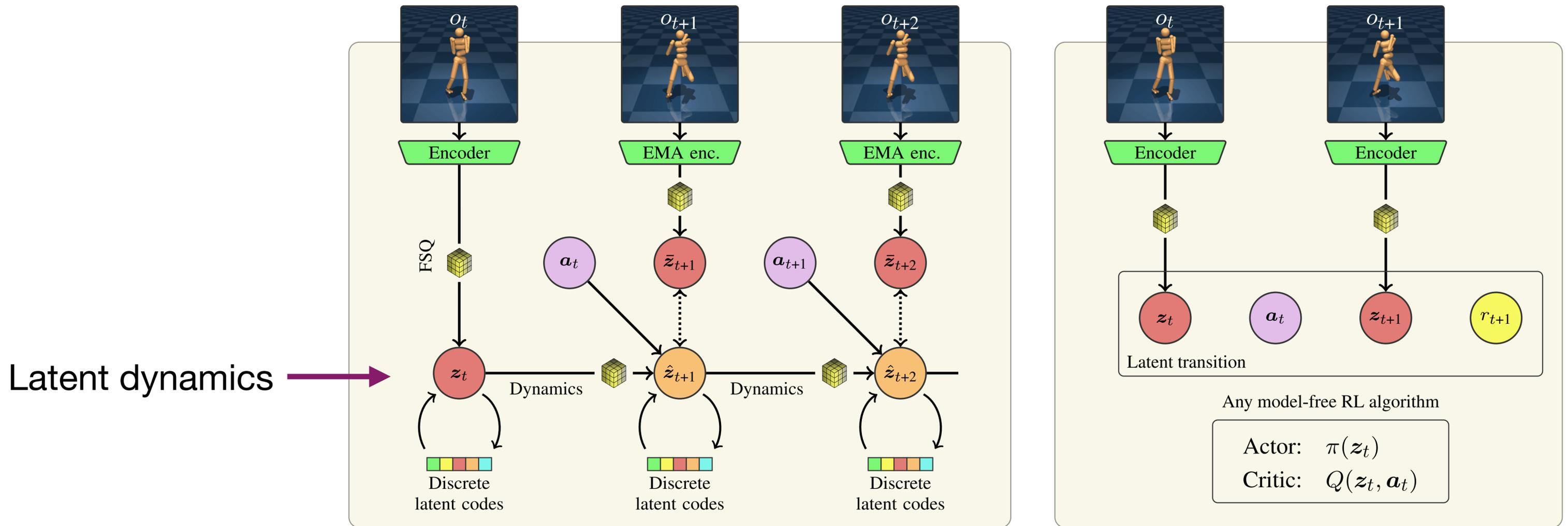
Quantization



# Case Study: iQRL

Dynamics Model but Not Model-based RL? 🤔

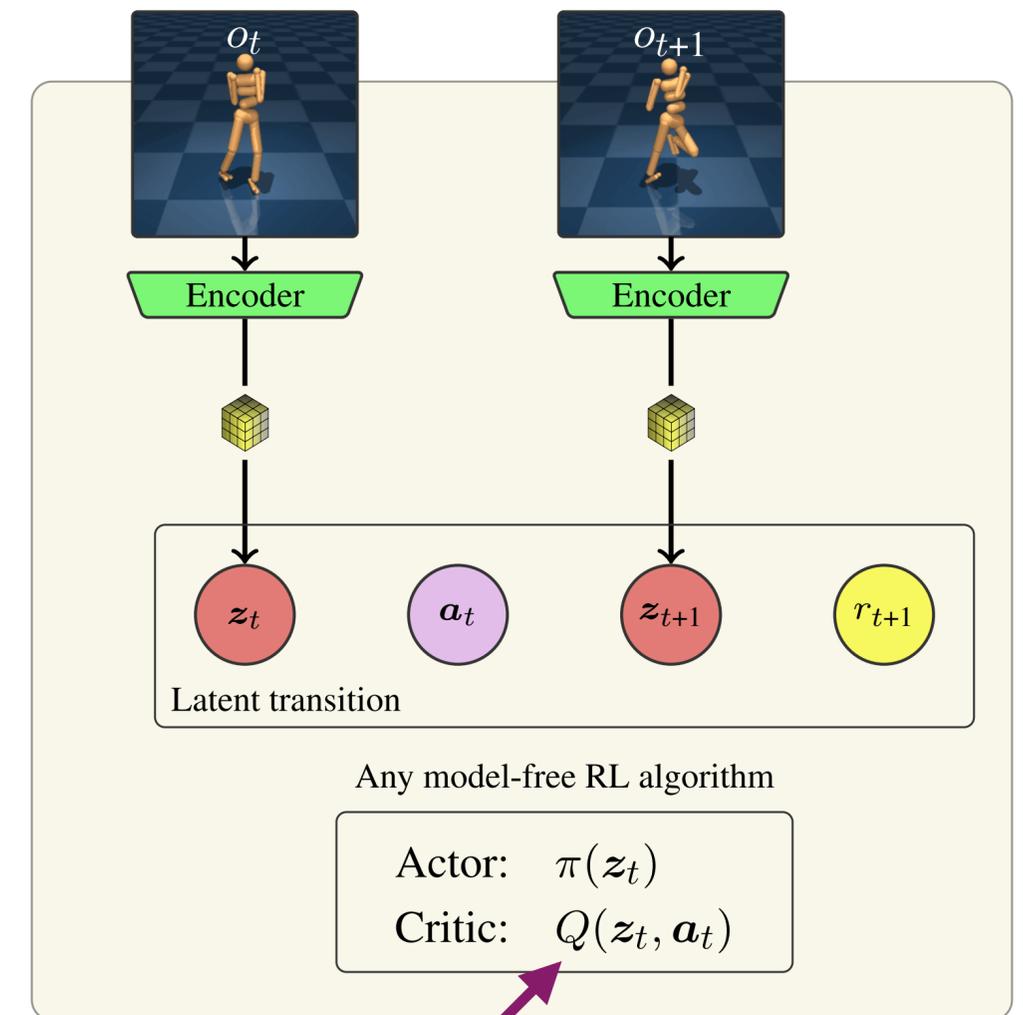
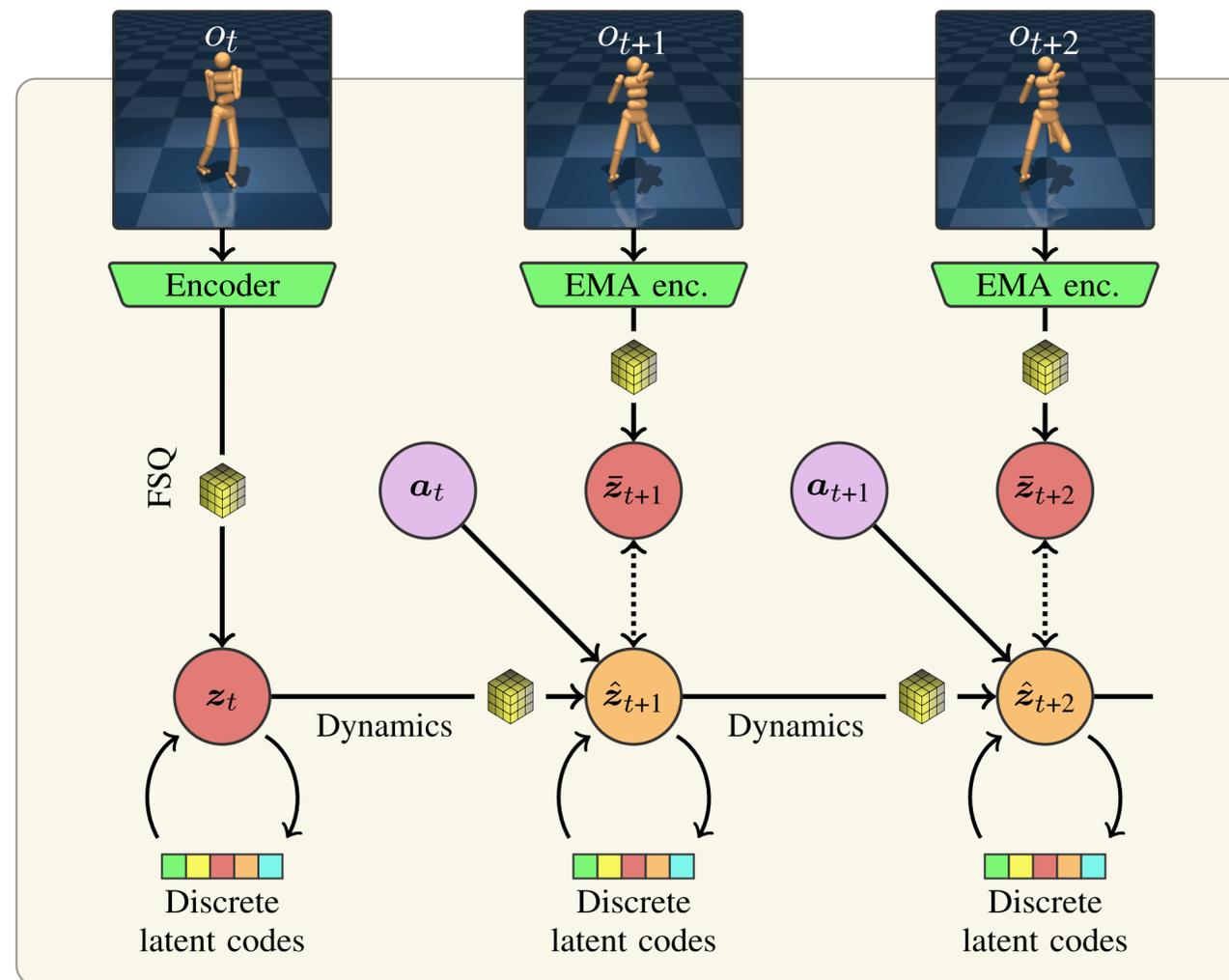
Use dynamics for representation learning



# Case Study: iQRL

Dynamics Model but Not Model-based RL? 🤔

Use dynamics for representation learning



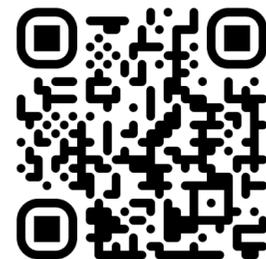
Model-free RL in latent space

# Outlook

# Outlook

**Email:** [aidan.scannell@aalto.fi](mailto:aidan.scannell@aalto.fi)

**Website:** [www.aidanscannell.com](http://www.aidanscannell.com)

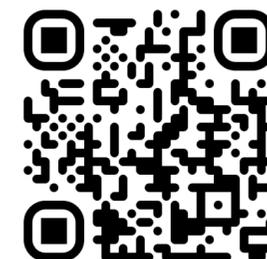


# Outlook

1. Model-based RL is a powerful tool

**Email:** [aidan.scannell@aalto.fi](mailto:aidan.scannell@aalto.fi)

**Website:** [www.aidanscannell.com](http://www.aidanscannell.com)

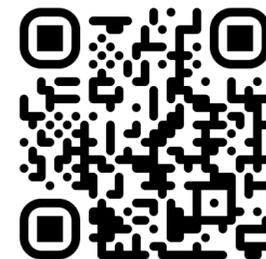


# Outlook

1. Model-based RL is a powerful tool
2. Leveraging predictive models improves sample efficiency

**Email:** [aidan.scannell@aalto.fi](mailto:aidan.scannell@aalto.fi)

**Website:** [www.aidanscannell.com](http://www.aidanscannell.com)



# Outlook

1. Model-based RL is a powerful tool
2. Leveraging predictive models improves sample efficiency
3. Lots more exciting work to be done

**Email:** [aidan.scannell@aalto.fi](mailto:aidan.scannell@aalto.fi)

**Website:** [www.aidanscannell.com](http://www.aidanscannell.com)

