

# Geometry of Neural Reinforcement Learning in Continuous State and Action Spaces

ICLR 2025 ORAL

Saket Tiwari (on job market for PostDoc)

Omer Gottesman

George Konidaris

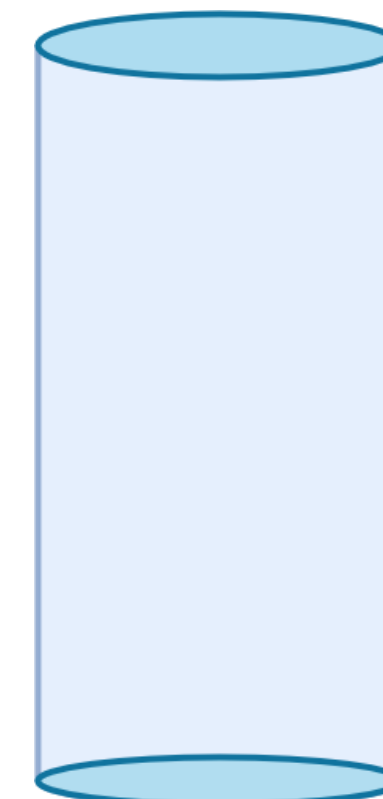
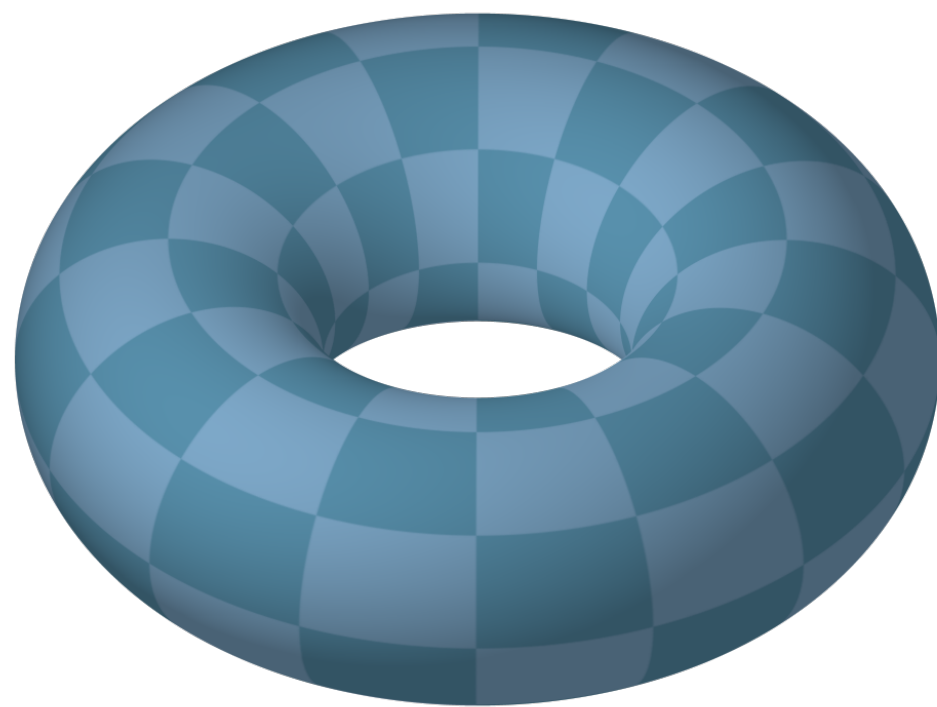
April, 2025

# Data Geometry and Manifold Hypothesis

Tenenbaum, et al 2000

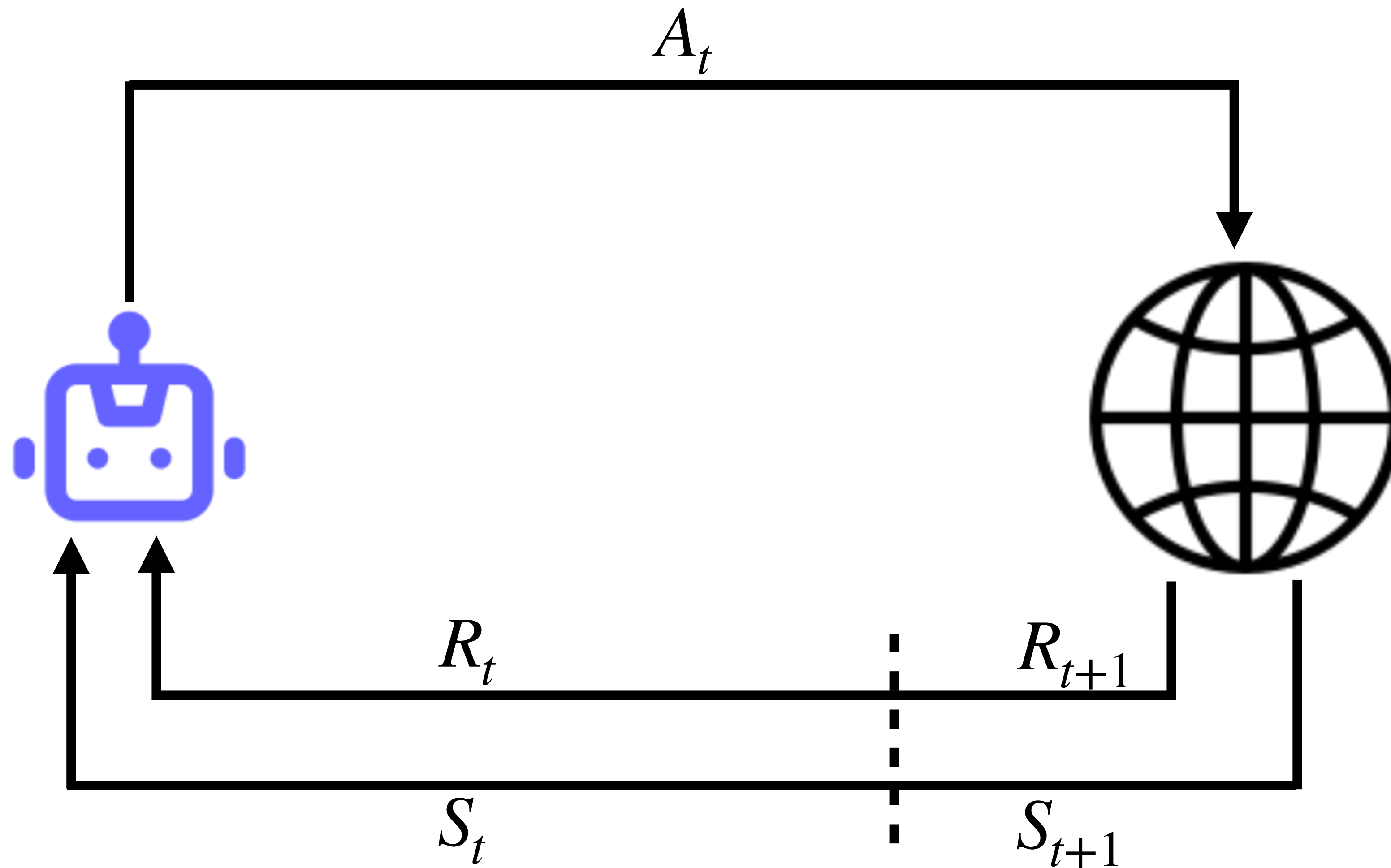
## The Manifold Hypothesis

Natural data lies close to lower-dimensional manifolds in its embedding space.

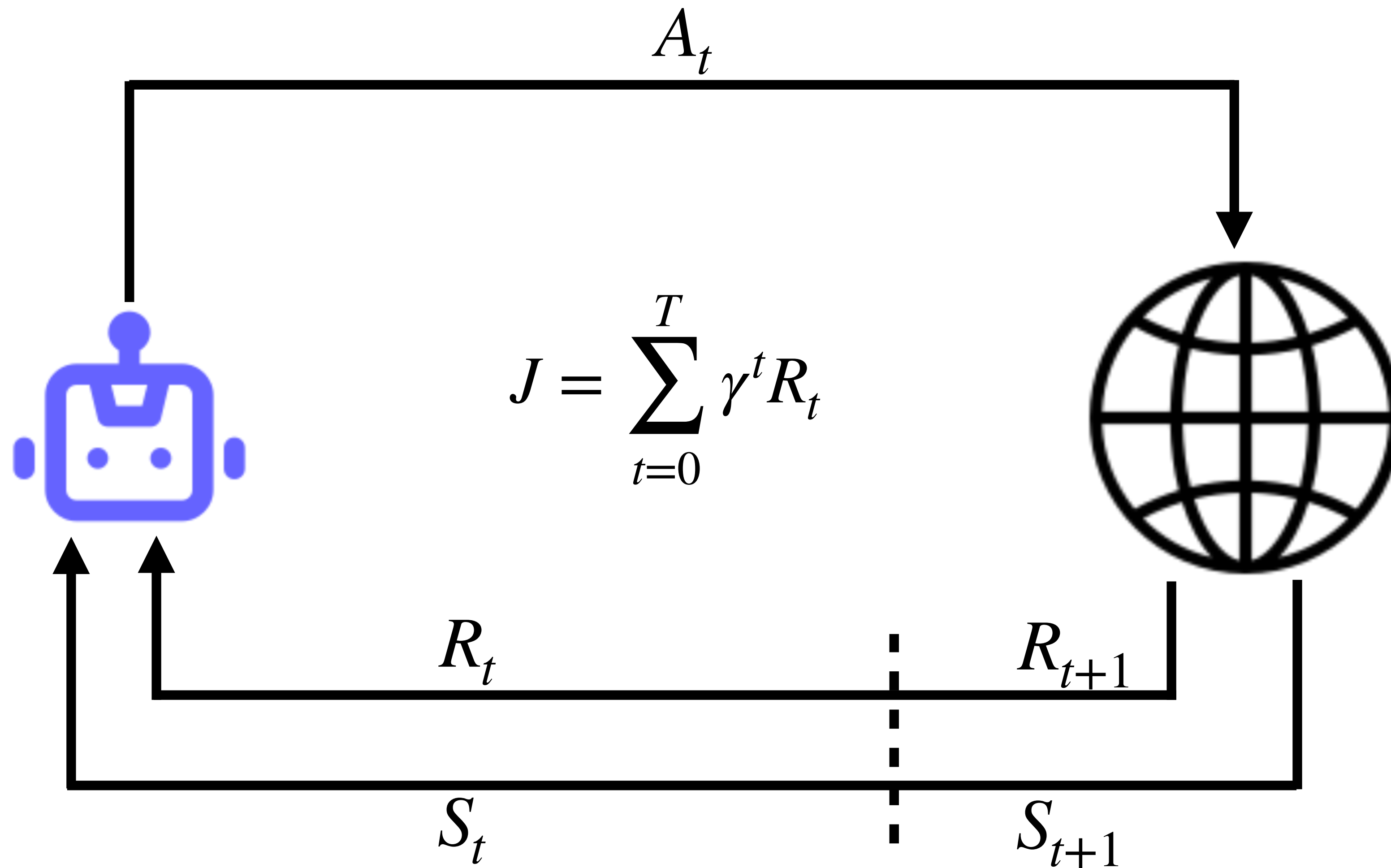


2D Manifolds embedded in 3D spaces: Torus, Sphere and Cylinder

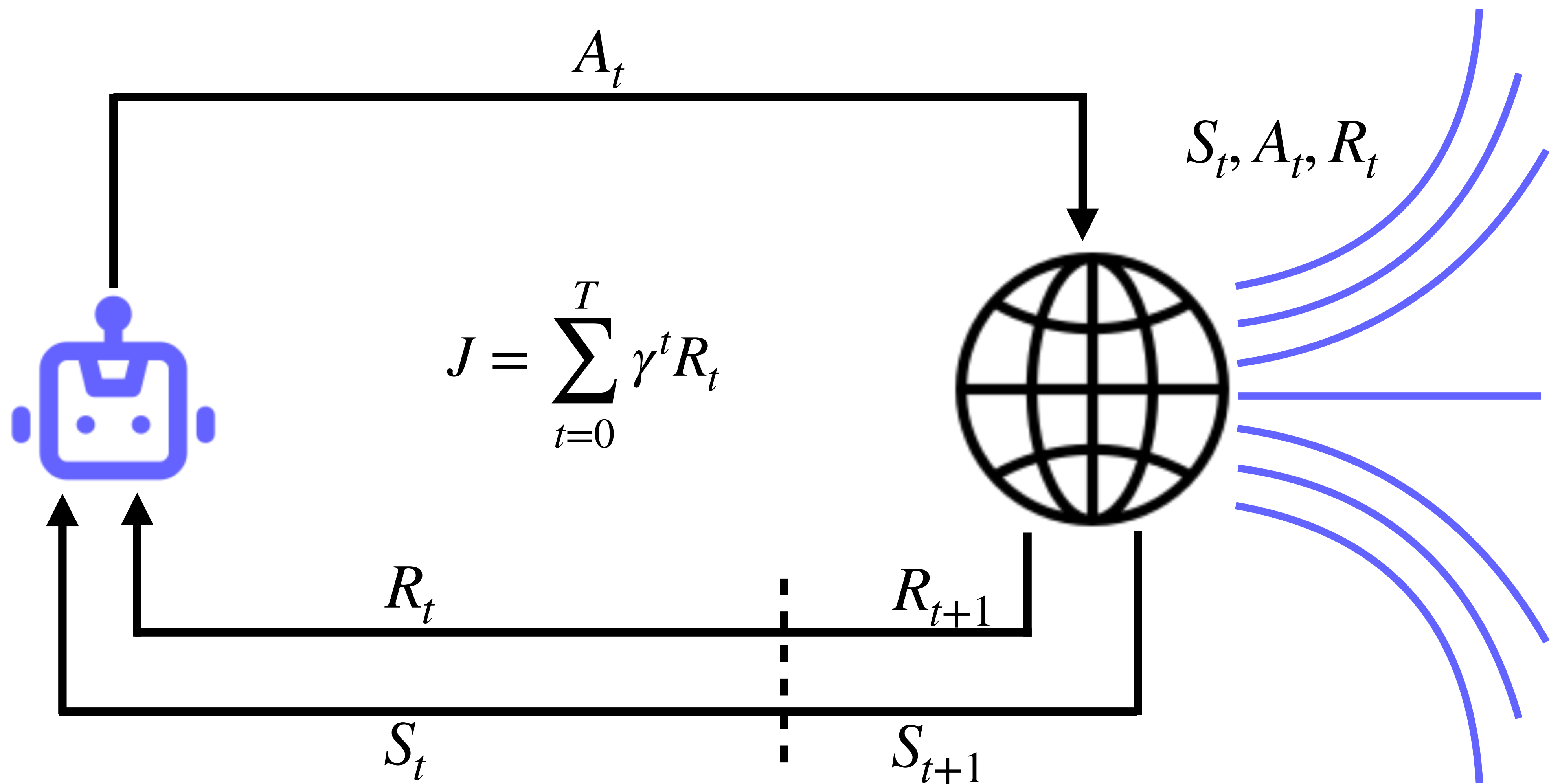
# Reinforcement Learning



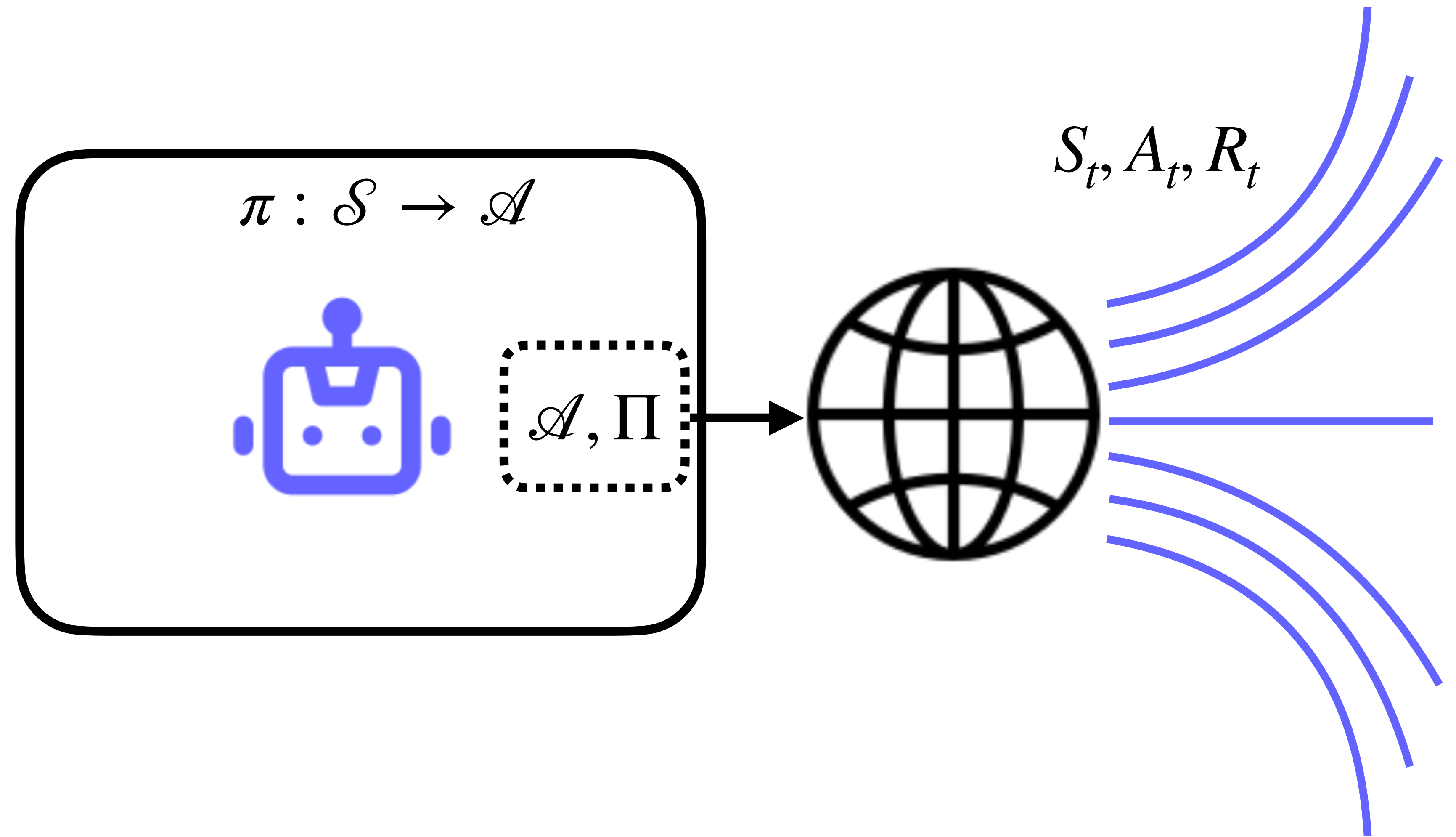
# Reinforcement Learning



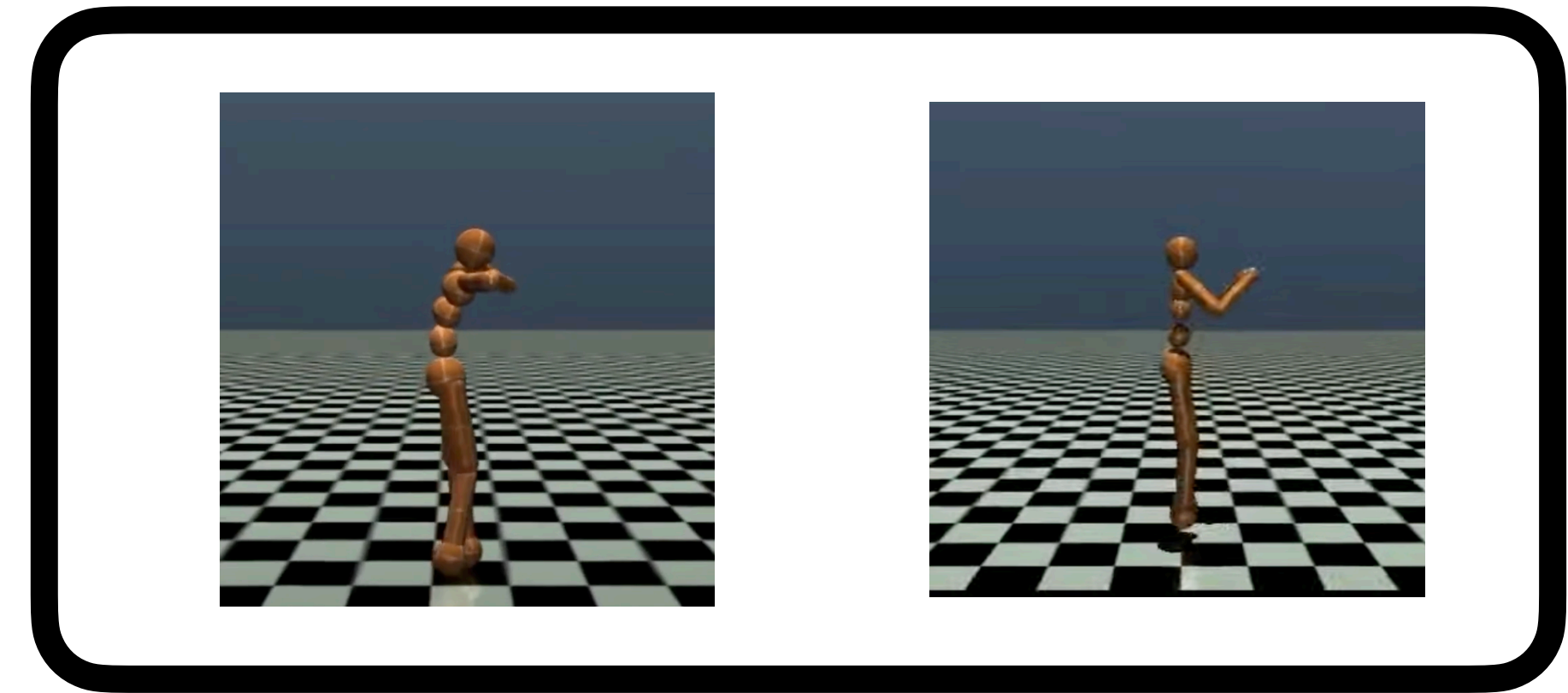
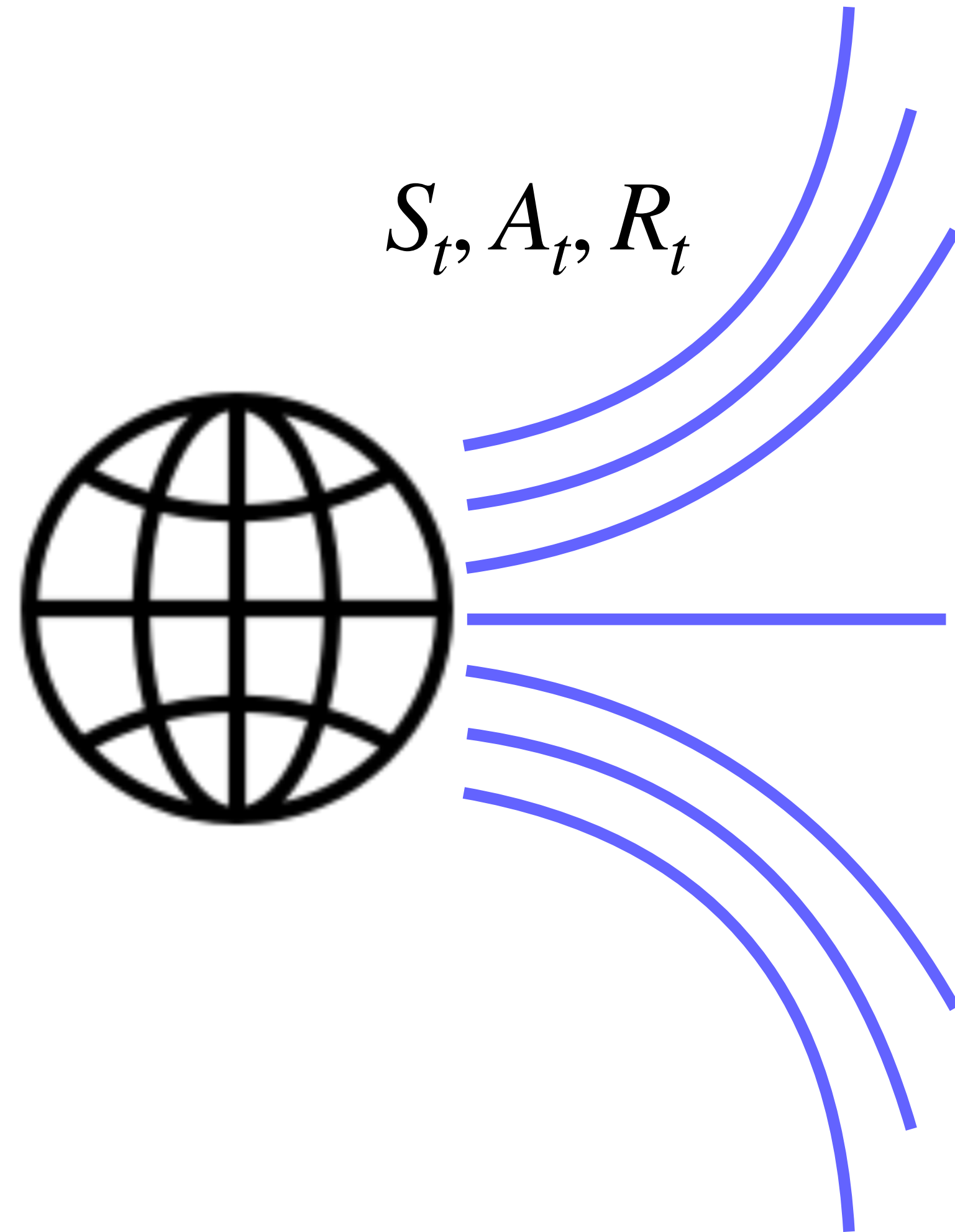
# Reinforcement Learning



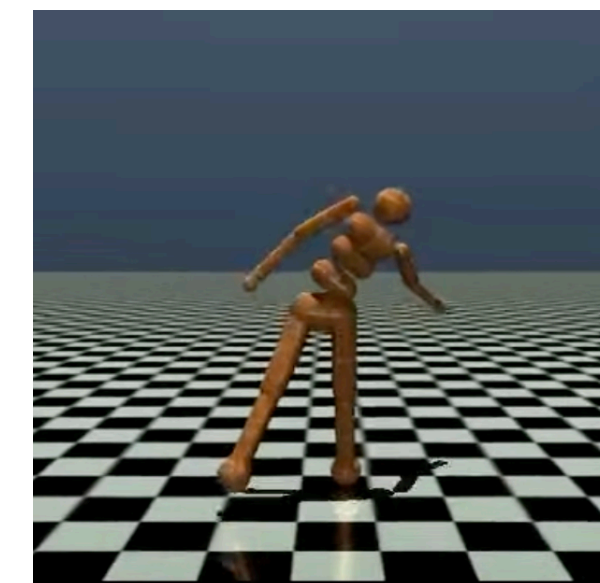
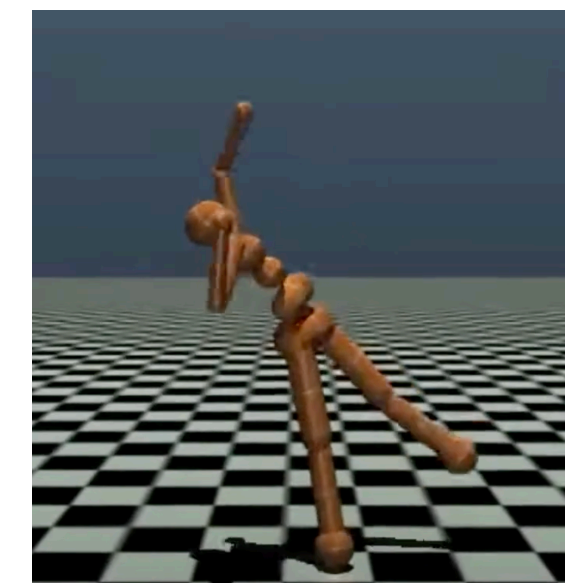
# Reinforcement Learning



# Reinforcement Learning



>>





# Reinforcement Learning: Policy Gradient Ascent

$$\pi(s; \theta_{\pi})$$

$$J(\theta) = \sum_{t=0}^T \gamma^t R_t$$

$$\theta_{\pi}^1 \leftarrow \theta_{\pi}^0 + \alpha \nabla_{\theta} J(\theta) \big|_{\theta=\theta_{\pi}^0}$$

$$\theta_{\pi}^0 \rightarrow \theta_{\pi}^1 \dots \rightarrow \theta_{\pi}^{\tau} \dots \rightarrow \theta_{\pi}^*$$



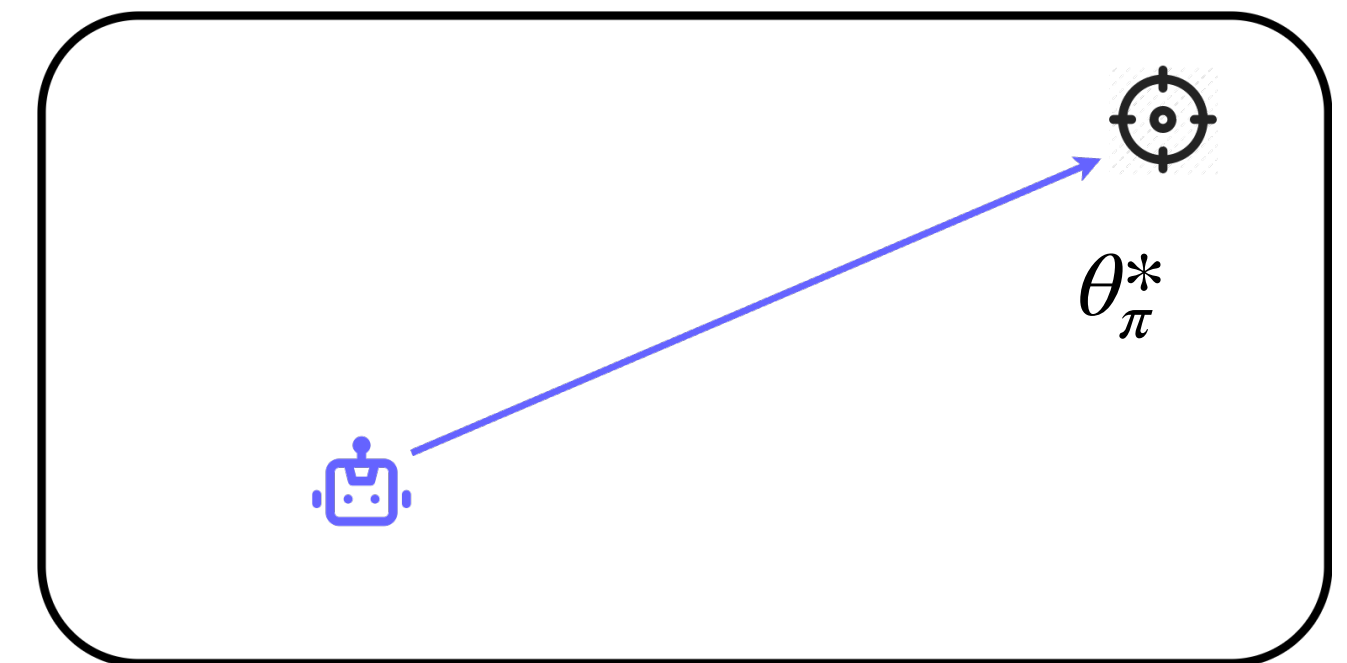
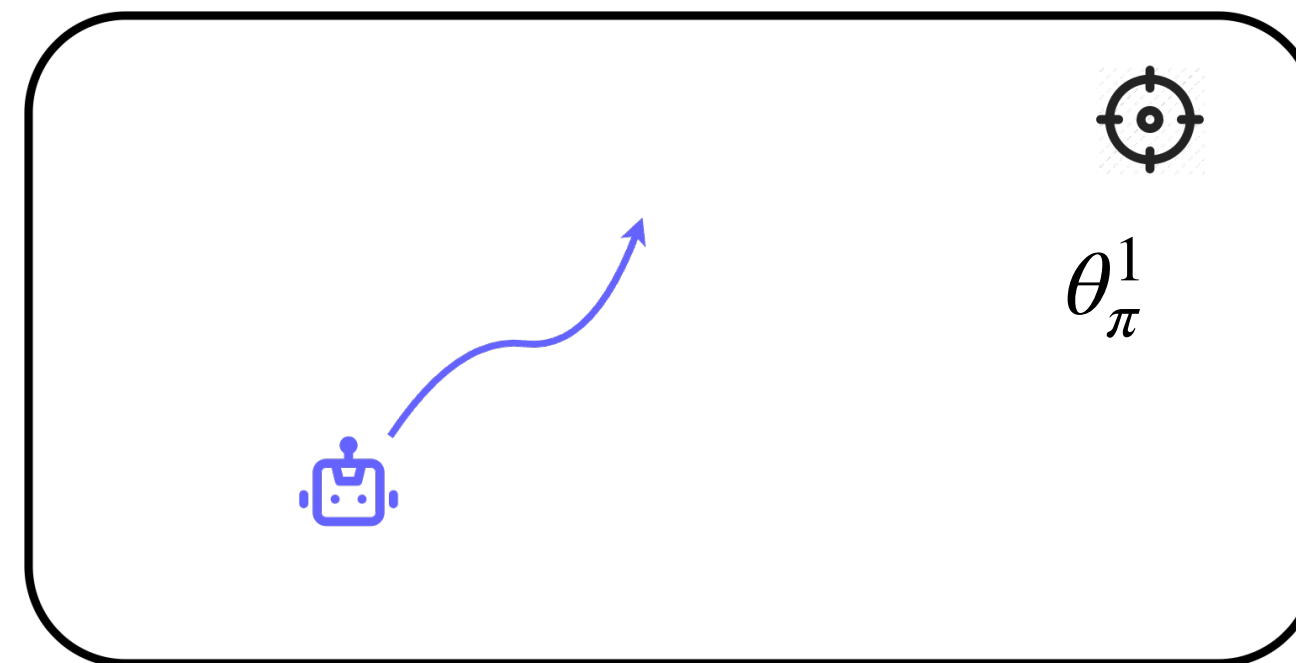
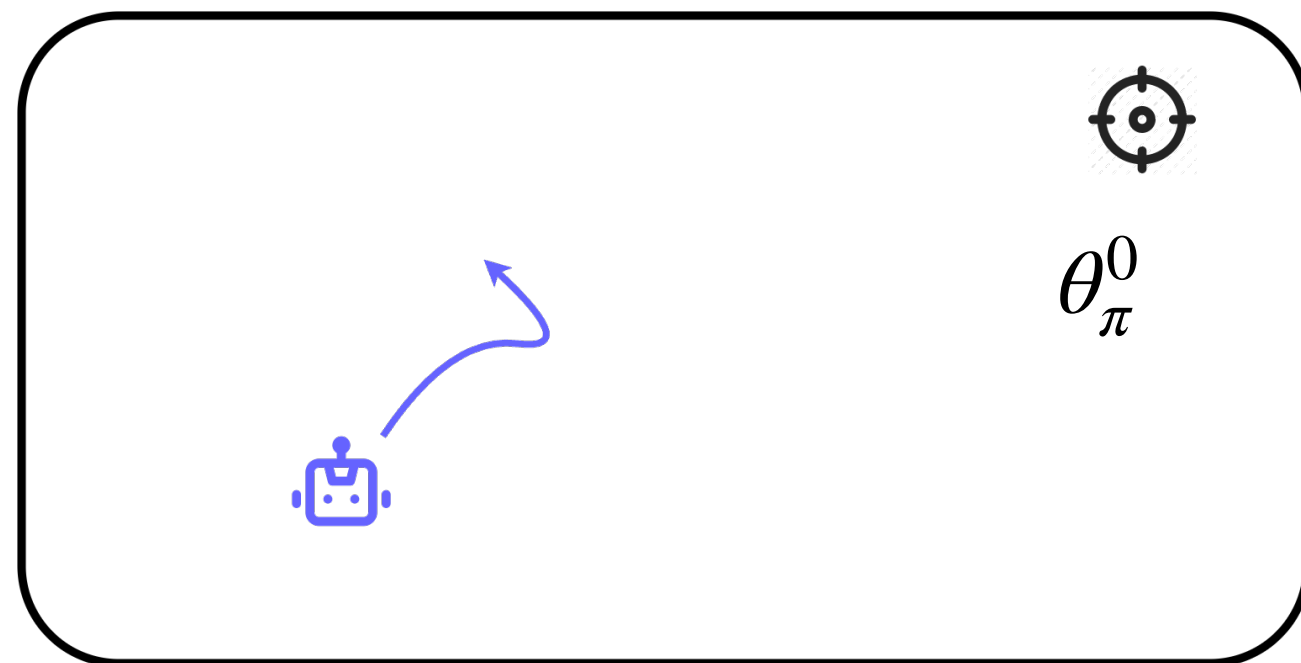
# Reinforcement Learning: Policy Gradient Ascent

$$\theta_{\pi}^0 \rightarrow \theta_{\pi}^1 \dots \rightarrow \theta_{\pi}^{\tau} \dots \rightarrow \theta_{\pi}^*$$



# Reinforcement Learning: Policy Gradient Ascent

$$\theta_{\pi}^0 \rightarrow \theta_{\pi}^1 \dots \rightarrow \theta_{\pi}^{\tau} \dots \rightarrow \theta_{\pi}^*$$

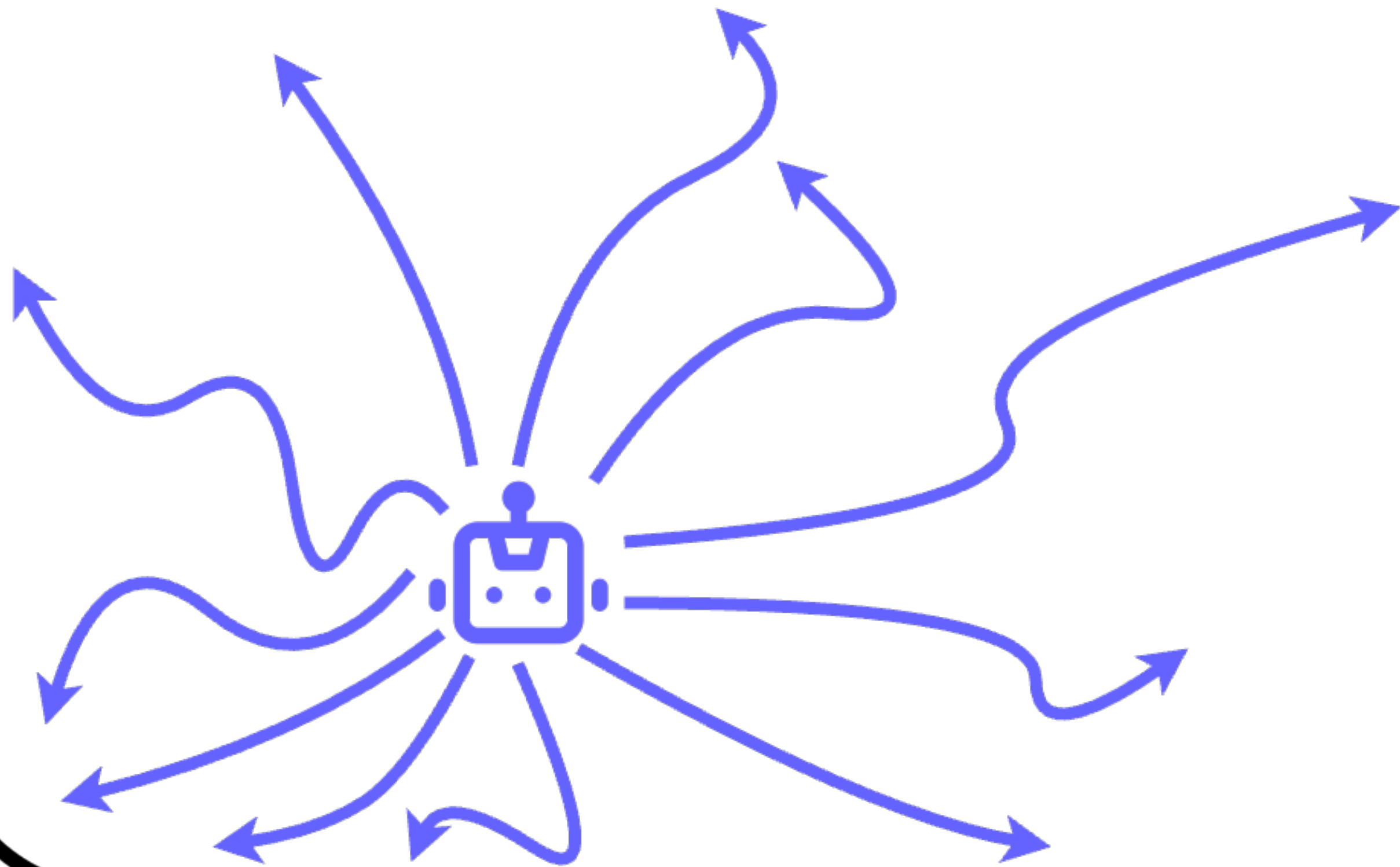
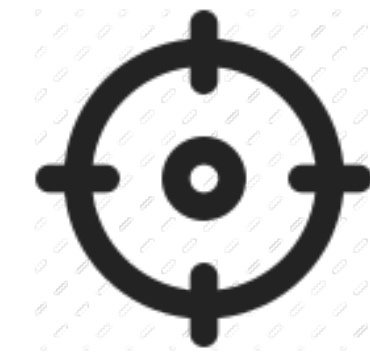


# Deep Reinforcement Learning

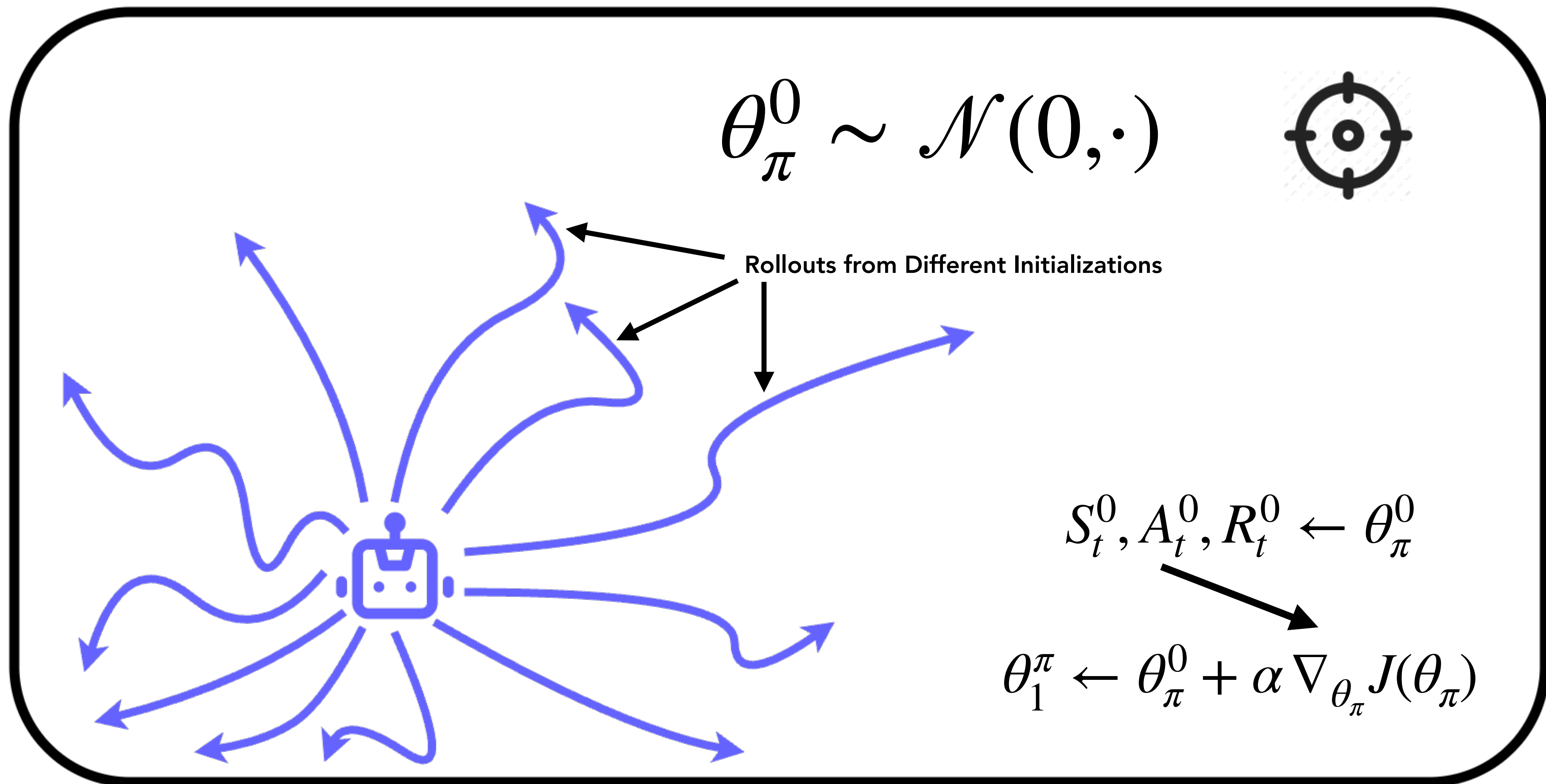
$$\theta_{\pi}^0 \sim \mathcal{N}(0, \cdot)$$

# Deep Reinforcement Learning

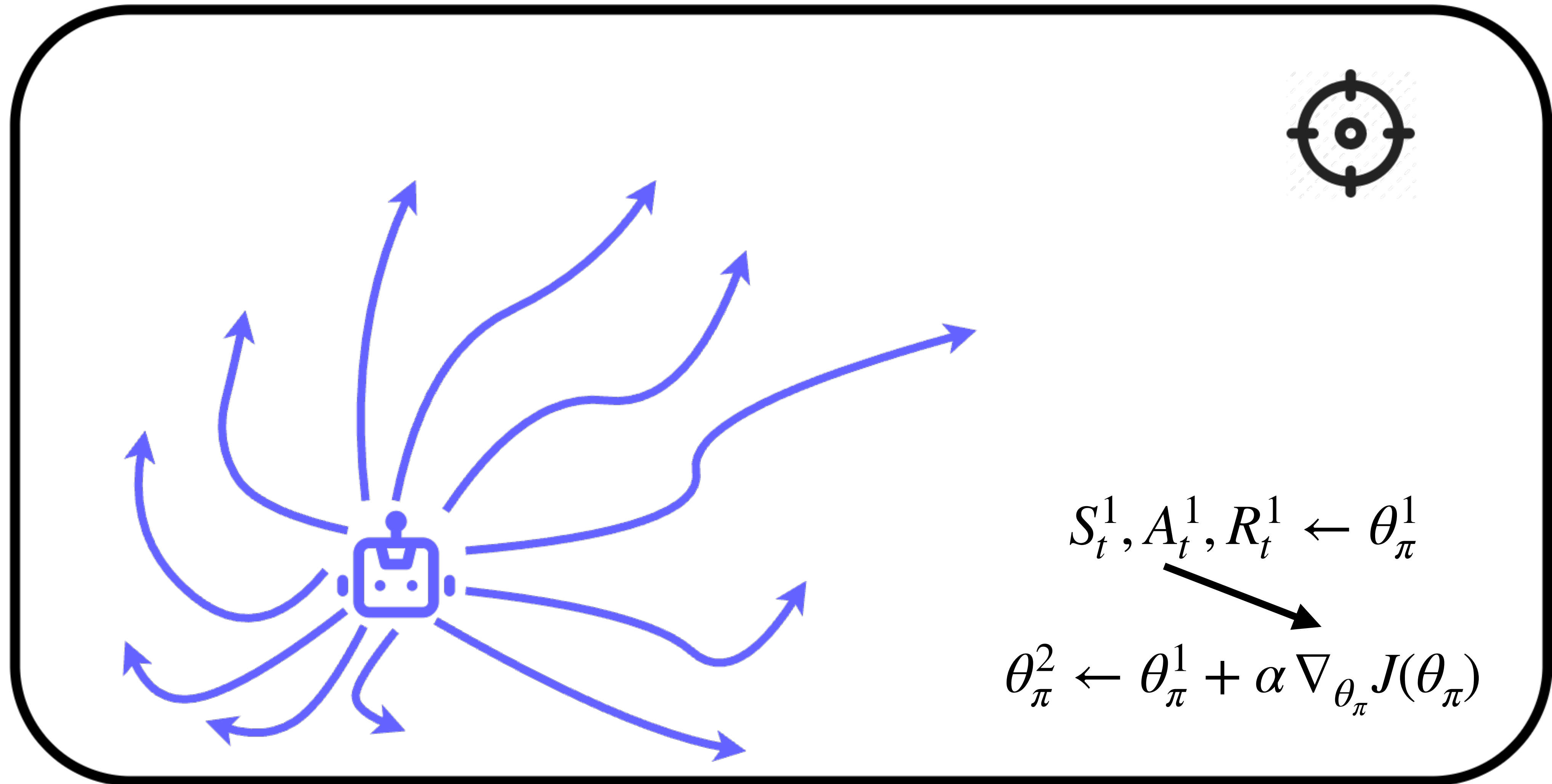
$$\theta_{\pi}^0 \sim \mathcal{N}(0, \cdot)$$



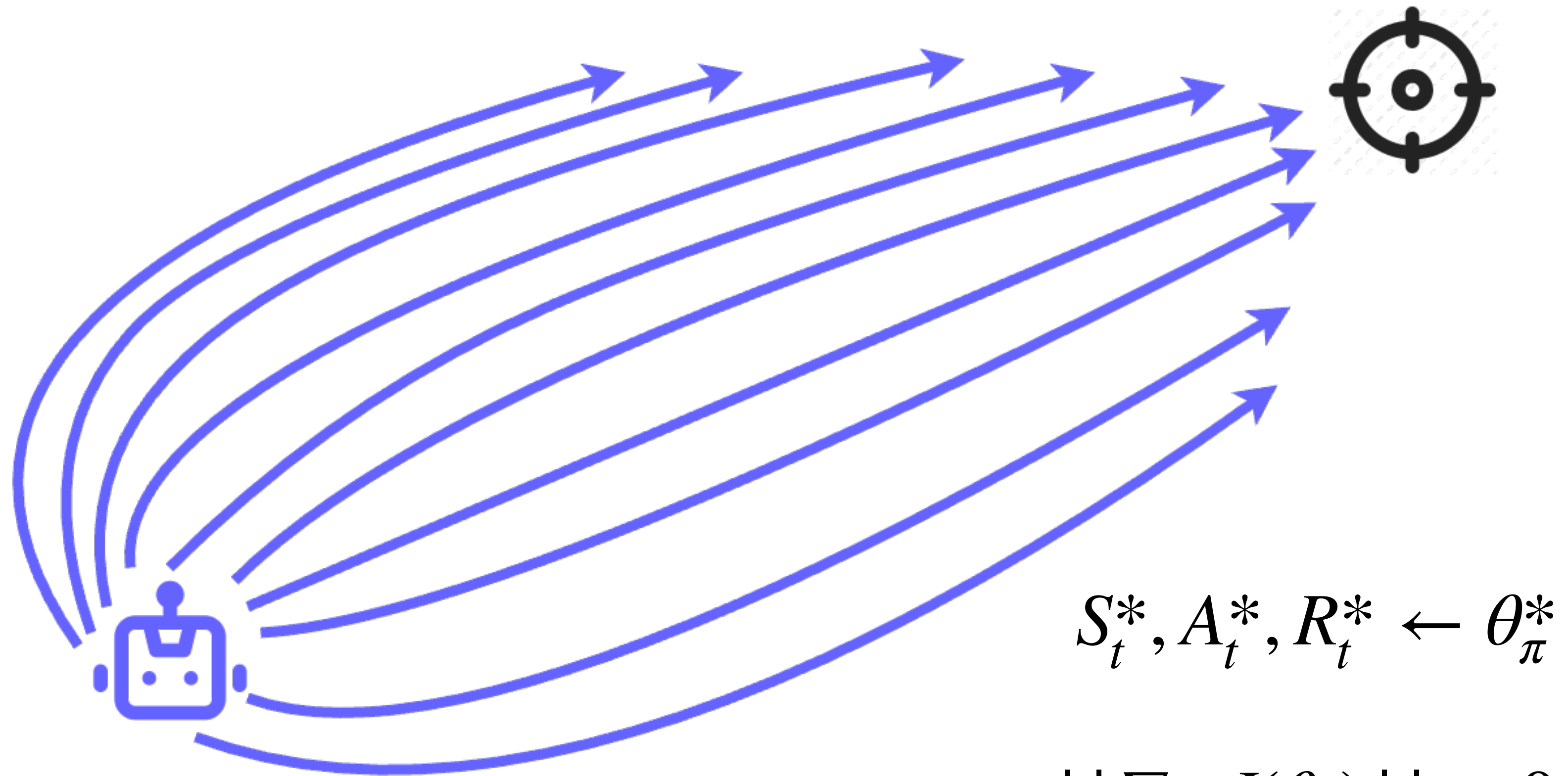
# Deep Reinforcement Learning



# Deep Reinforcement Learning



# Deep Reinforcement Learning



$$S_t^*, A_t^*, R_t^* \leftarrow \theta_\pi^*$$

$$||\nabla_{\theta_\pi} J(\theta_\pi)|| \approx 0$$



# Deep Reinforcement Learning

Distribution of parameters

$$\mathcal{N}(0, \cdot) = \mathbf{W}_{\pi}^0 \rightarrow \mathbf{W}_{\pi}^1 \dots \rightarrow \mathbf{W}_{\pi}^{\tau} \dots \rightarrow \mathbf{W}_{\pi}^*$$

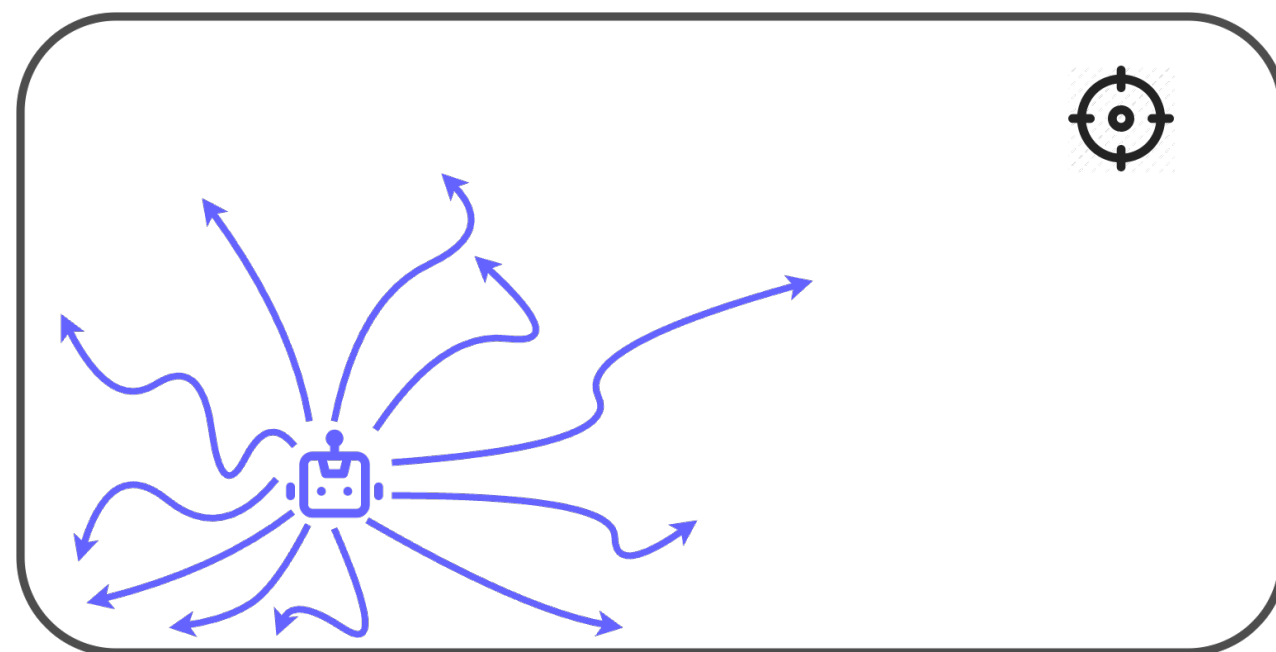
# Deep Reinforcement Learning

Distribution of parameters

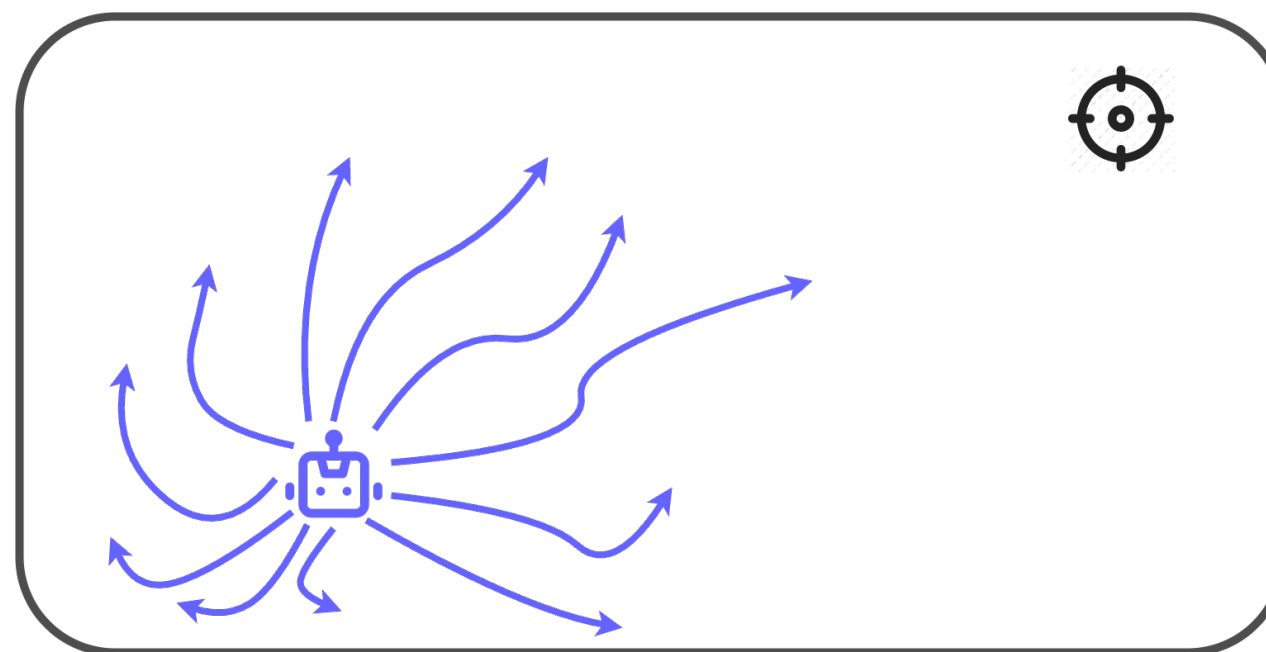
$$\mathcal{N}(0, \cdot) = \mathbf{W}_{\pi}^0 \rightarrow \mathbf{W}_{\pi}^1 \dots \rightarrow \mathbf{W}_{\pi}^{\tau} \dots \rightarrow \mathbf{W}_{\pi}^*$$

$$S_t^{\tau}$$

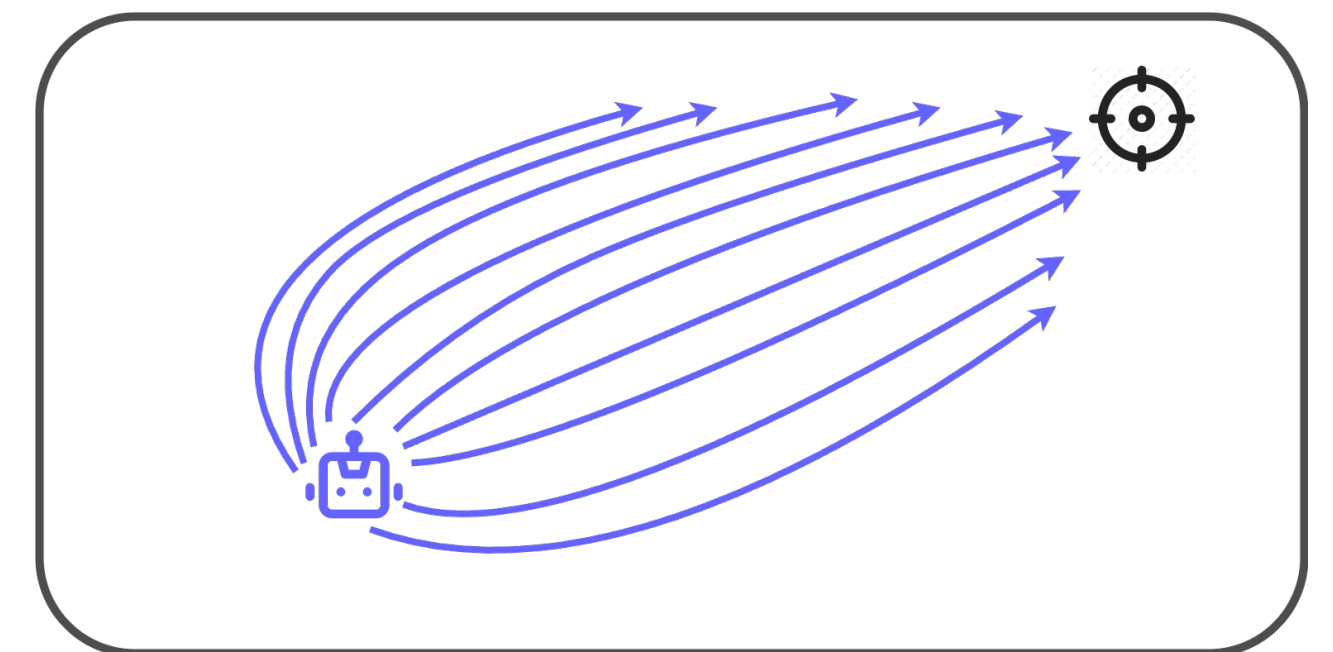
$\tau = 0$



$\tau = 1$



$\tau \rightarrow \infty$



# Low-Dimensional Data in RL for Deterministic Environments

## Theoretical Setting

- Deterministic transitions
- Continuous states, actions, time
- Neural networks in the infinite width limit

# Continuous time

**Why?** Gradient Descent vs Gradient Flow

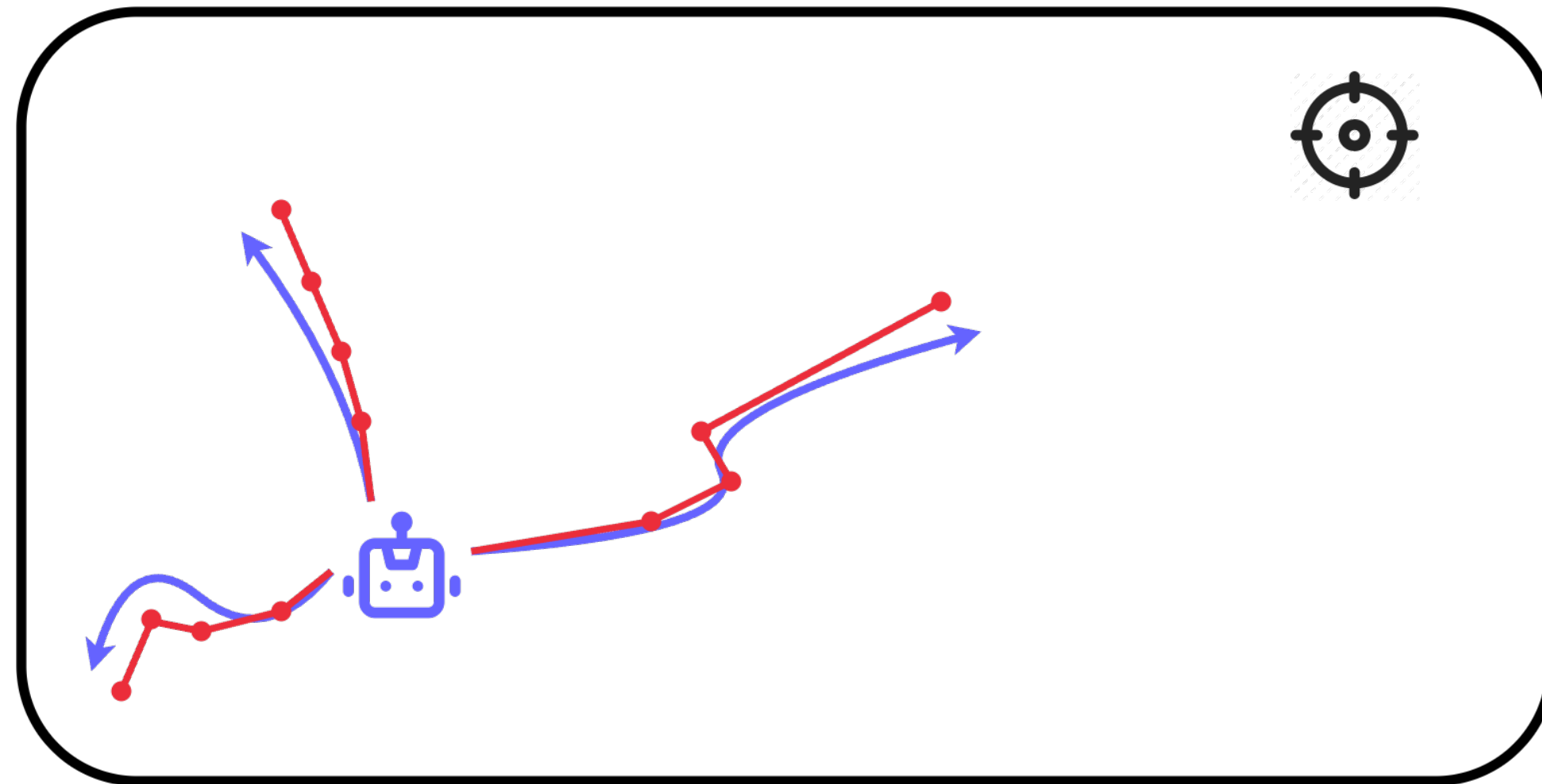
$$x_{t+1} \leftarrow x_t - \alpha \nabla_x L(x)$$

VS

$$\frac{dx(t)}{dt} = - \nabla L(x(t))$$

# Continuous time

## Continuous Time vs Discrete Time RL



**Baird, 1994**

**Doya, 2000**

**Munos, 2006**

**Wang et al. 2020**


# Continuous state, action, time

## Continuous Time Deterministic MDP

$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, (g, H), r, \gamma, T, s_0 \rangle$$

# Continuous state, action, time

## Continuous Time Deterministic MDP

$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, (g, H), r, \gamma, T, s_0 \rangle$$


$\mathcal{S} \subseteq \mathbb{R}^{d_s}$      $\mathcal{A} \subseteq \mathbb{R}^{d_a}$      $r(s_t)$      $s(0) = s_0$



# Continuous state, action, time

## Continuous Time Deterministic MDP

$$\begin{array}{c} \text{Transition Dynamics} \\ \downarrow \\ \mathcal{M} = \langle \mathcal{S}, \mathcal{A}, (g, H), r, \gamma, T, s_0 \rangle \end{array}$$
$$\frac{ds_t}{dt} = g(s_t) + H(s_t)\pi(s_t)$$
$$\begin{array}{ccc} \downarrow & & \downarrow \\ \mathbb{R}^{d_s} & & \mathbb{R}^{d_s \times d_a} \end{array}$$

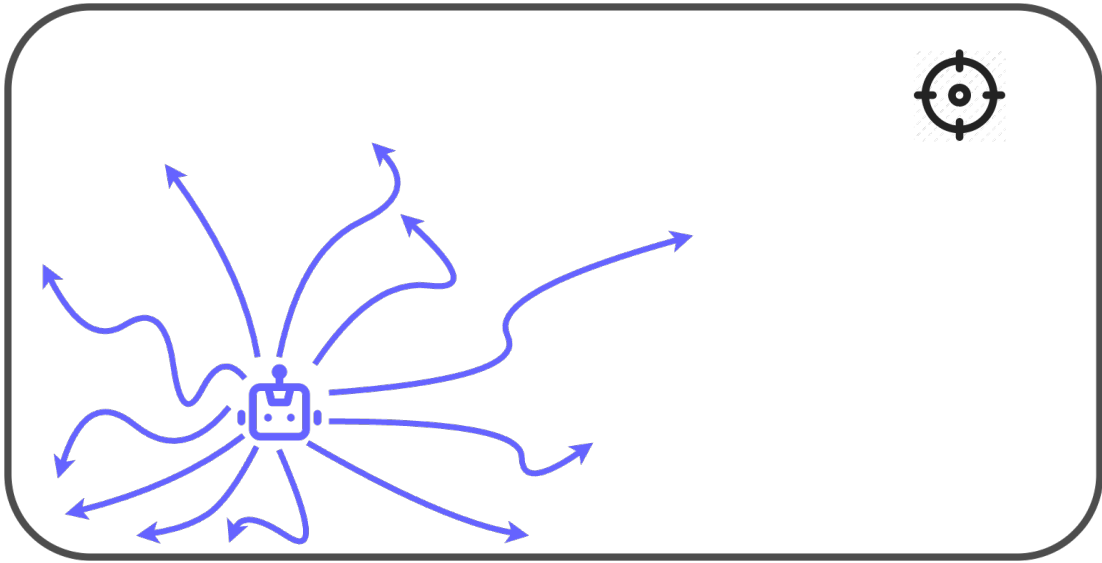
# Low-Dimensional Data in RL for Deterministic Environments

$$\frac{ds_t}{dt} = g(s_t) + H(s_t) \Phi(s_t, \mathbf{W}_\pi^0) \mathbf{W}_\pi^\tau$$

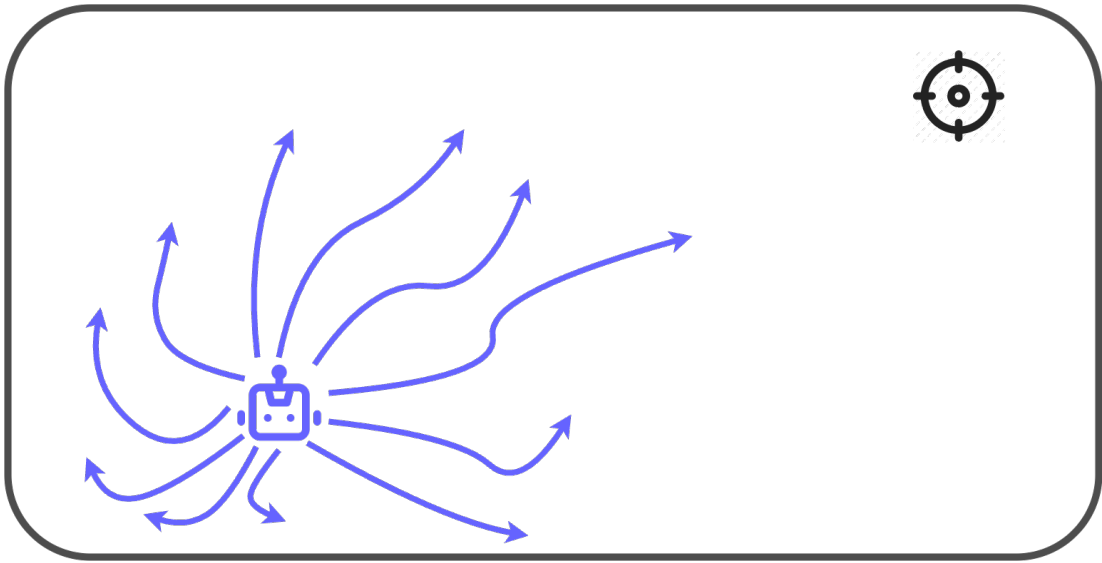
Infinite Width Neural Network

$S_t^\tau$

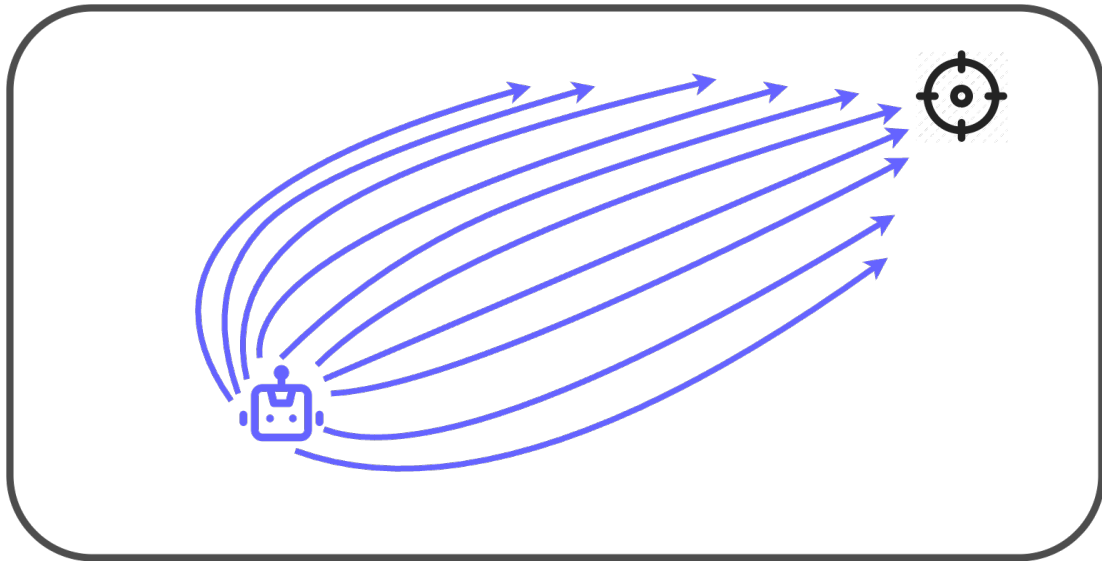
$\tau = 0$



$\tau = 1$



$\tau \rightarrow \infty$



# Low-Dimensional Data in RL for Deterministic Environments

Main result:

For  $t \in (0, \delta)$ , starting from a fixed state  $s$  :

$$S_t^\tau(n) \xrightarrow[\text{In probability}]{\text{width } n \rightarrow \infty} \hat{S}_t^\tau + O(\delta^3)$$

# Low-Dimensional Data in RL for Deterministic Environments

Main result:

For  $t \in (0, \delta)$ , starting from a fixed state  $s$  :

$$S_t^\tau(n) \xrightarrow[\text{In probability}]{\text{width } n \rightarrow \infty} \hat{S}_t^\tau + O(\delta^3)$$

$$\Pr(\text{distance}(\hat{S}_t^\tau, M) \geq D) \leq e^{-\mathcal{R}_\tau D}$$

$$\dim(M) \leq 2d_a + 1$$

# Low-Dimensional Data in RL for Deterministic Environments

Informal main result:

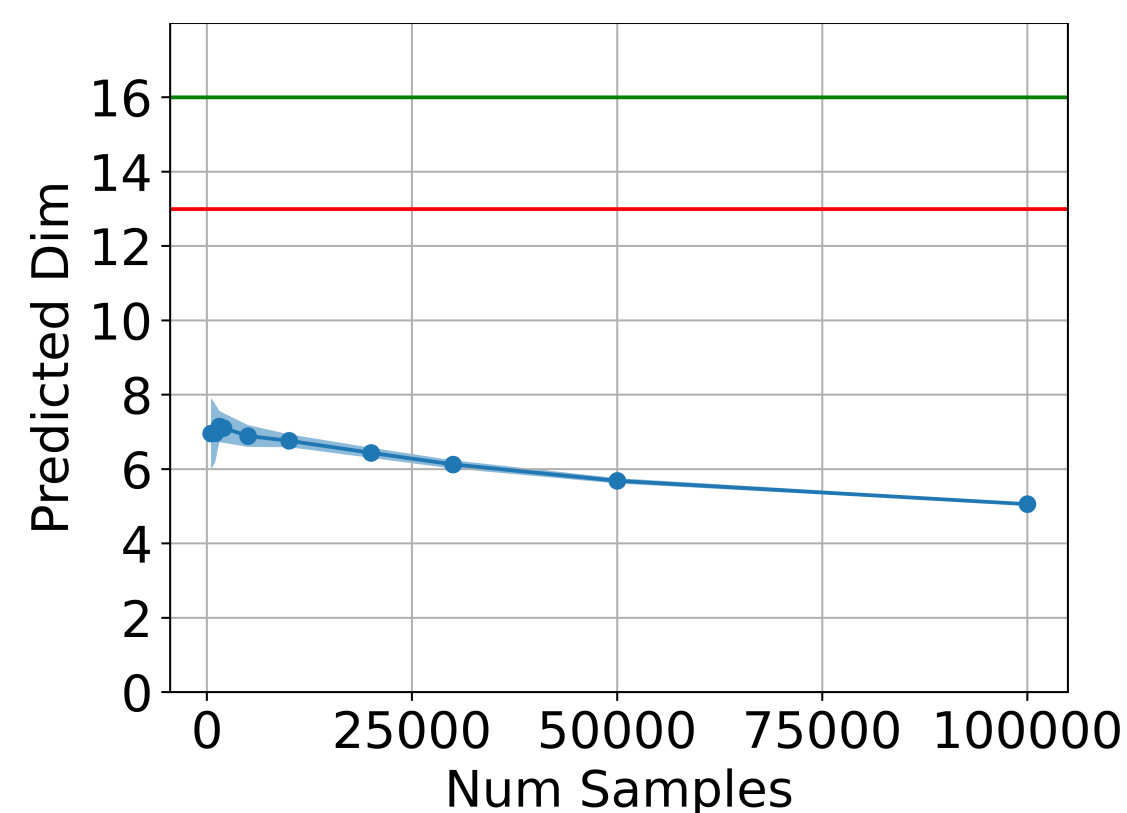
For a “well behaved” policy gradient algorithm with wide neural network policy and small learning rates: the agent’s states are concentrated around a low dimensional manifold.

# Empirical evidence: dimensionality estimation

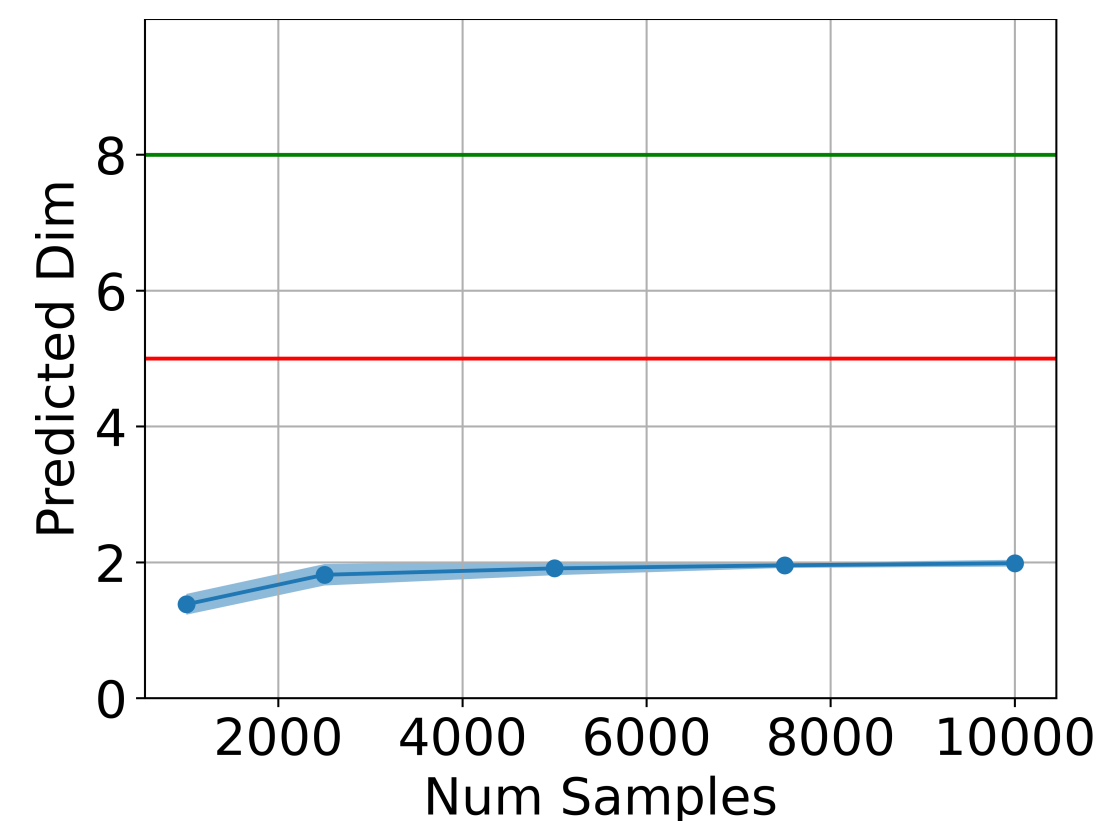
Facco et al 2017: empirical estimates

Green  $= d_s$ , Red  $= 2d_a + 1$

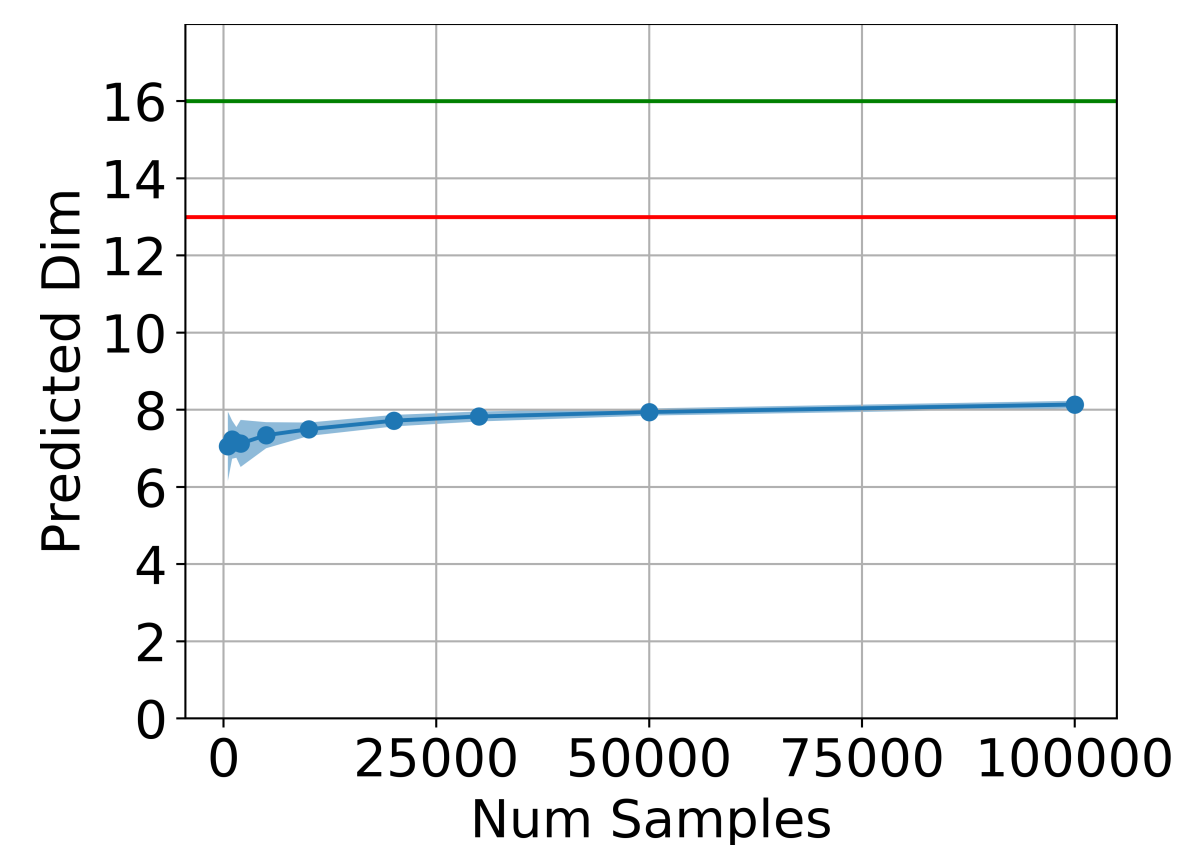
Walker



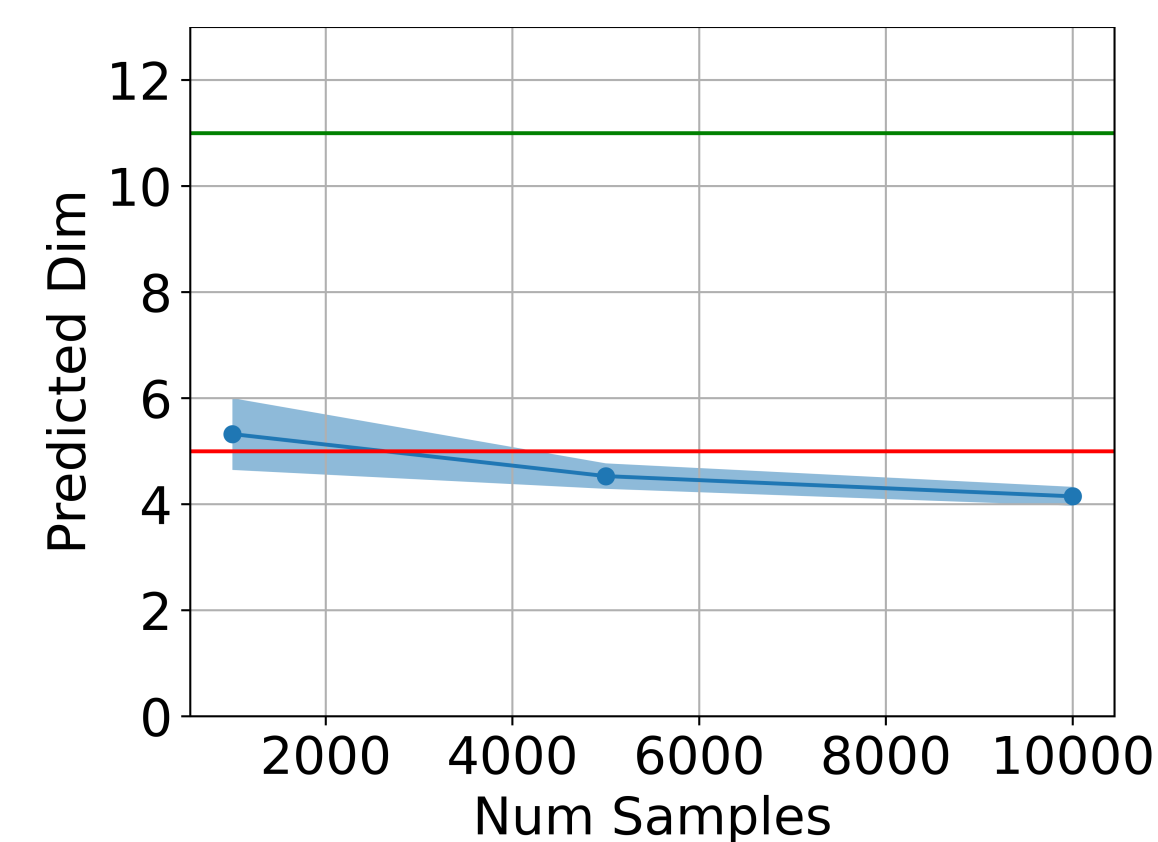
Pendulum



Cheetah



Reacher

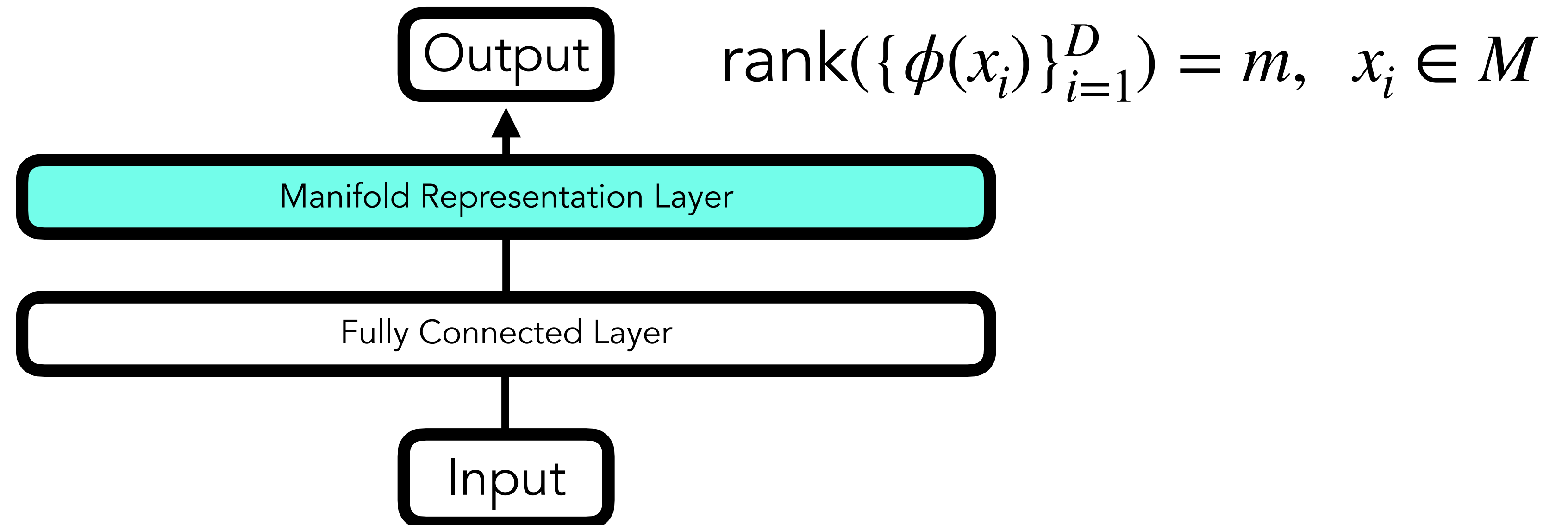


# Application: Exploit the Underlying Structure

## Manifold Representation Layer

Can we learn low-dimensional representations as theory suggests?

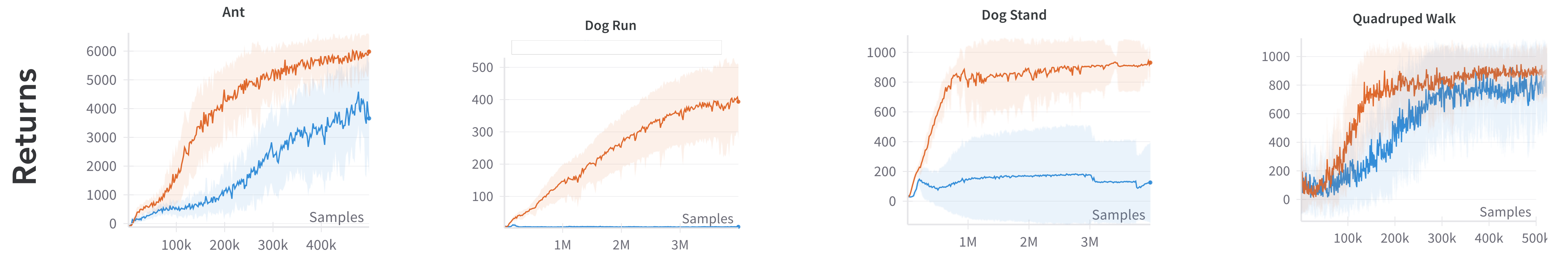
Yu et al. 2023, White box transformer





# Manifold Representation Layer

**Red** = Ours, **Blue** = Baseline



There is a **low-dimensional structure** to data in reinforcement learning and neural networks are well suited to **exploit this structure**.

# Hire me for postdoc!



[saket.ai](https://saket.ai)