

TDDE78 Reinforcement Learning

Value-Based and Policy Based RL

Amath Sow

PhD student ReaL Lab

Reinforcement Learning: Concepts



AGENT

Agent: takes actions.

Reinforcement Learning: Concepts



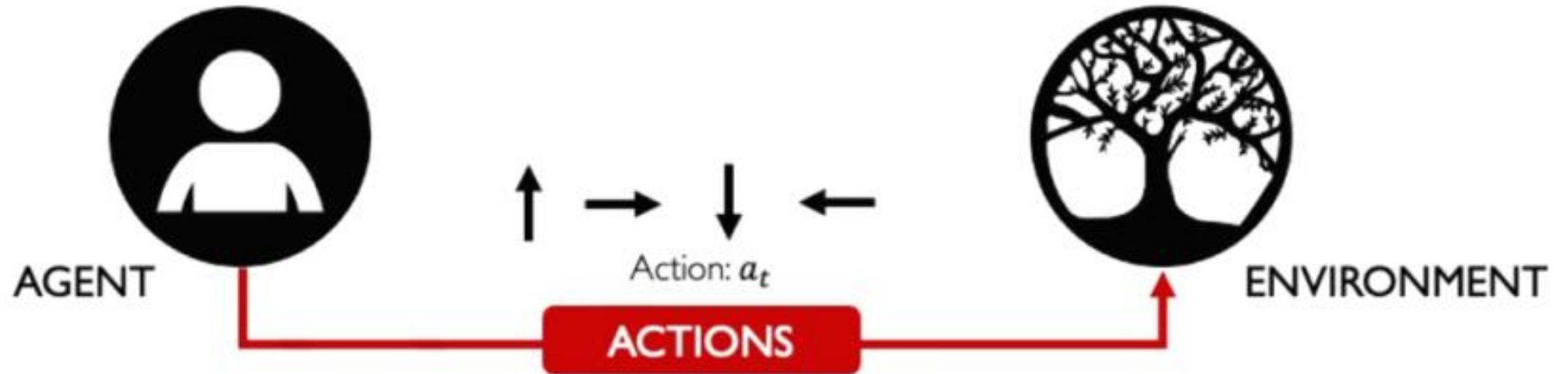
AGENT



ENVIRONMENT

Environment: the world in which the agent exists and operates.

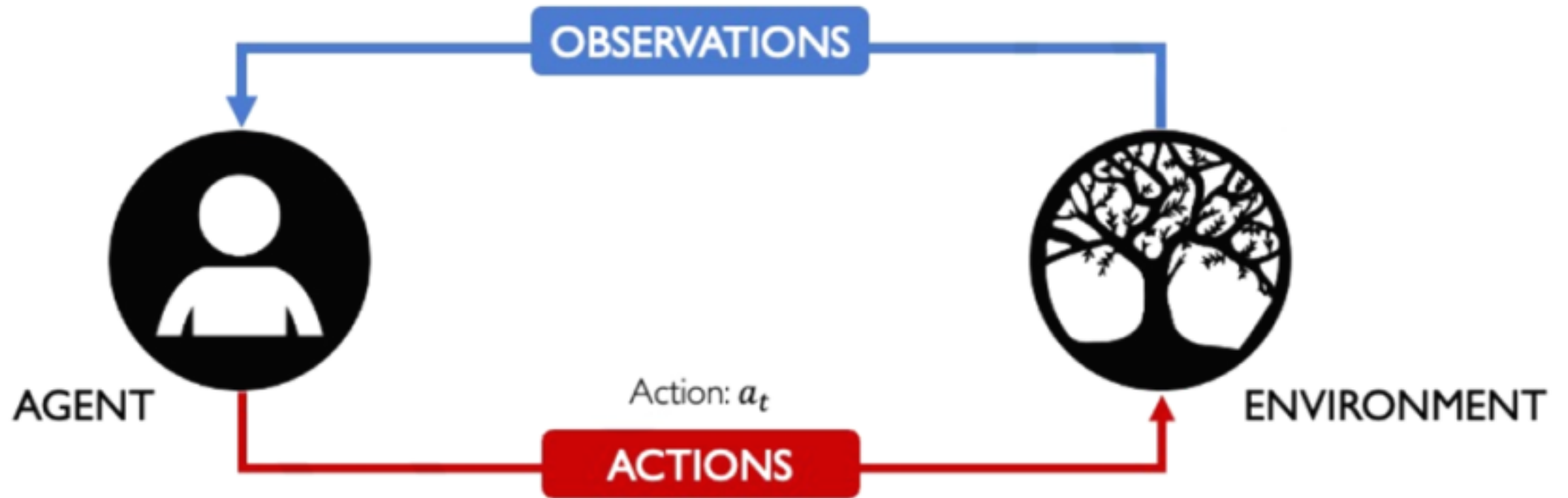
Reinforcement Learning: Concepts



Action: a move the agent can make in the environment.

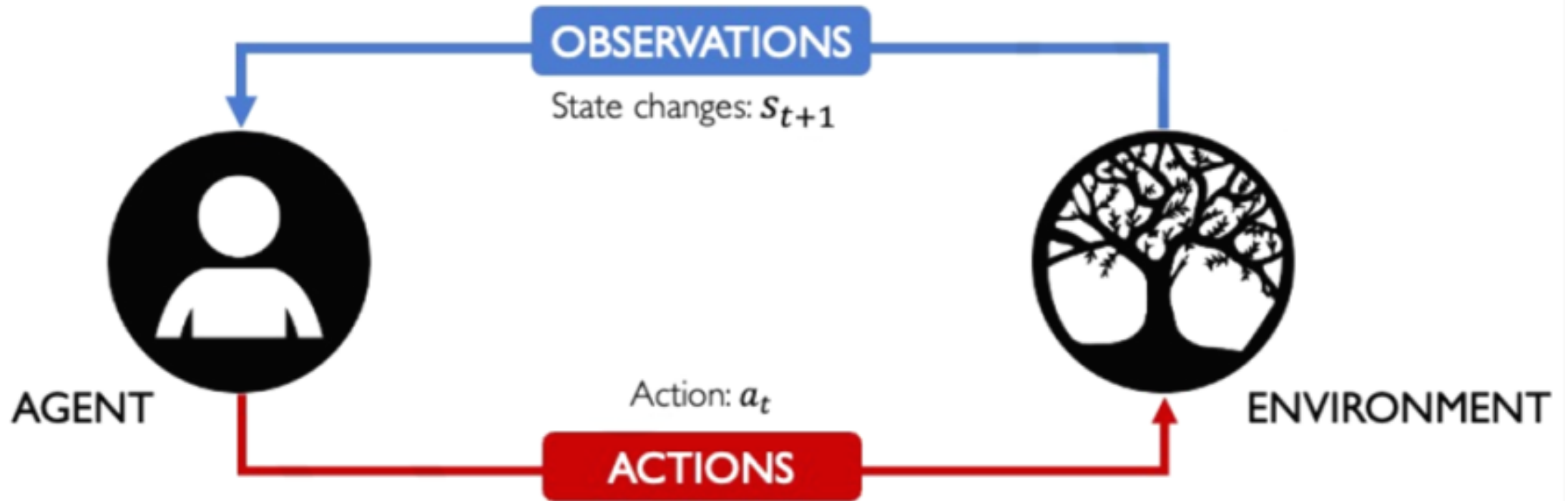
Action space A : the set of possible actions an agent can make in the environment

Reinforcement Learning: Concepts



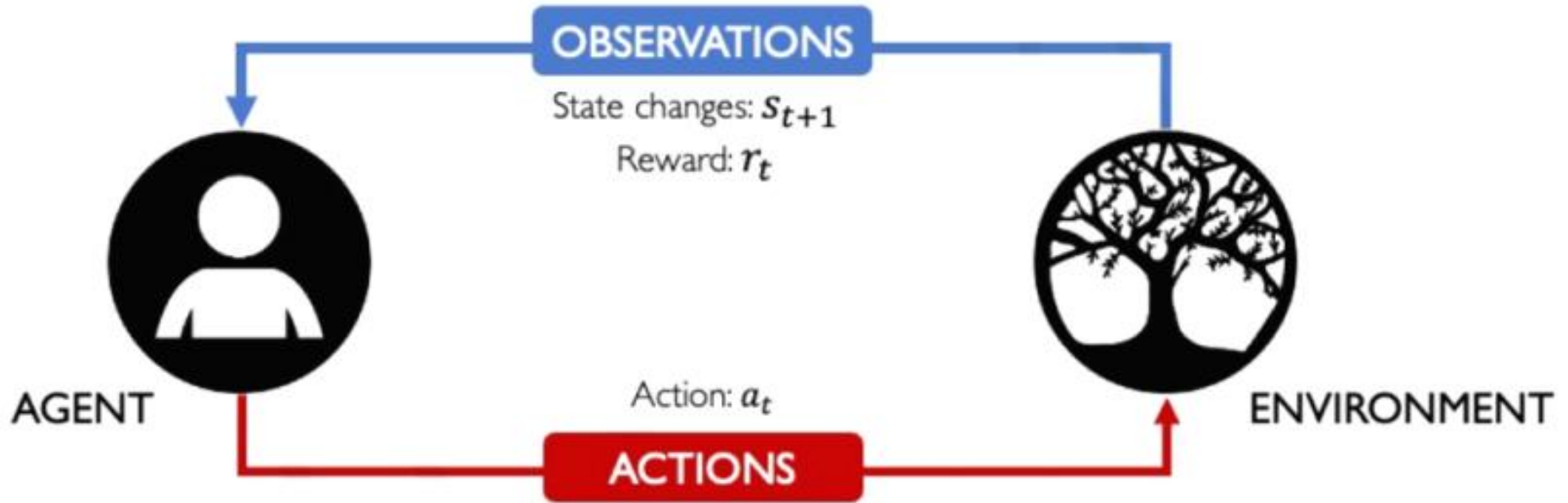
Observations: of the environment after taking actions.

Reinforcement Learning: Concepts



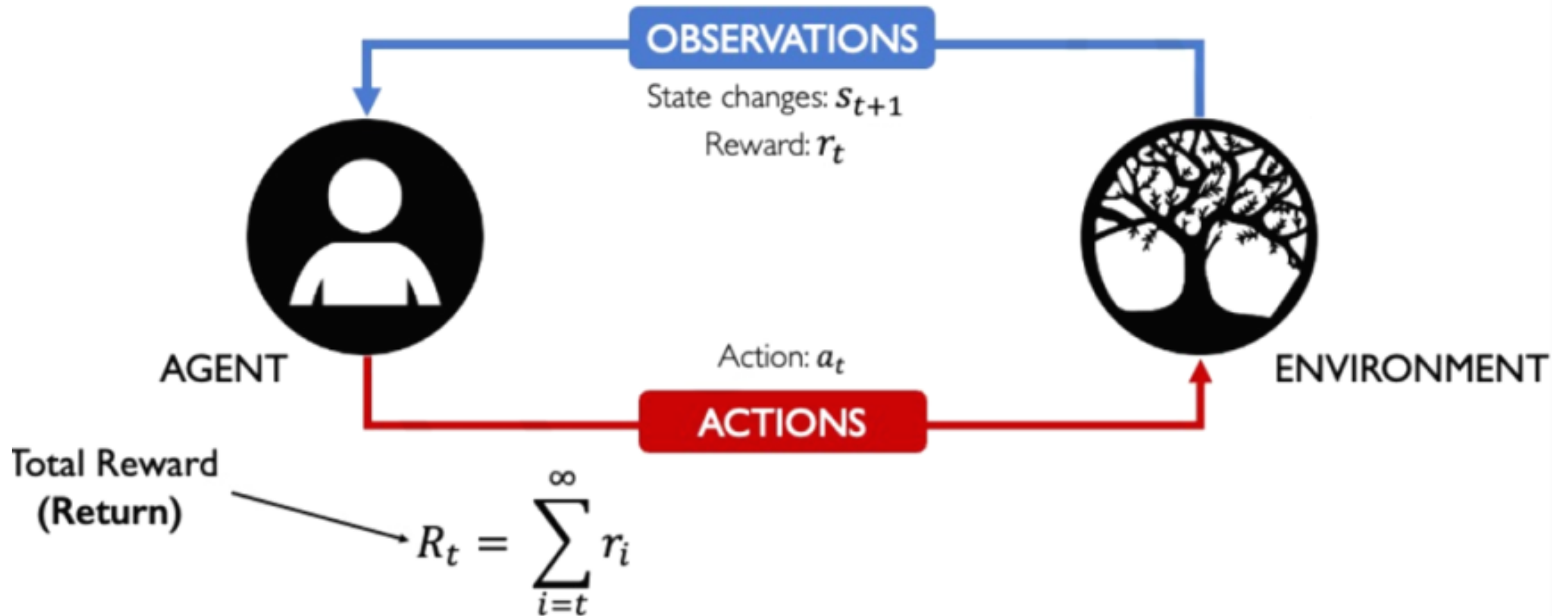
State: a situation which the agent perceives.

Reinforcement Learning: Concepts

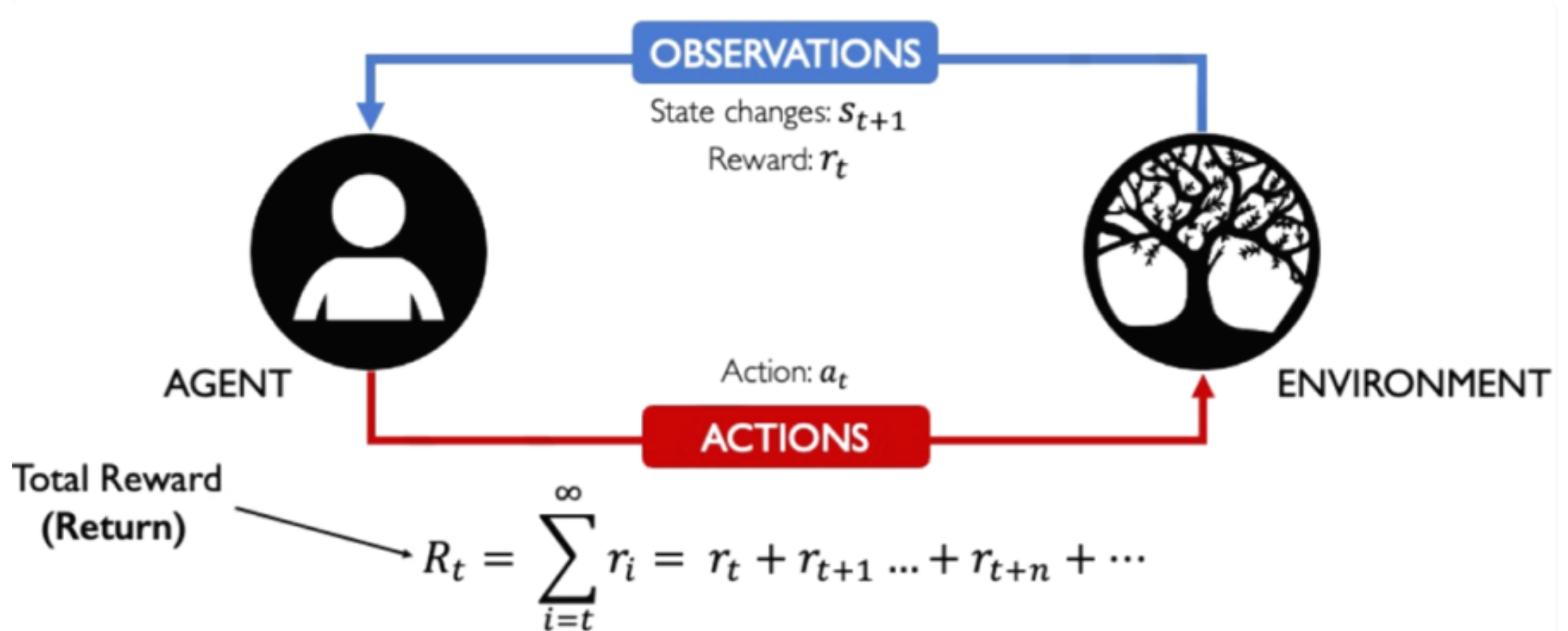


Reward: feedback that measures the success or failure of the agent's action.

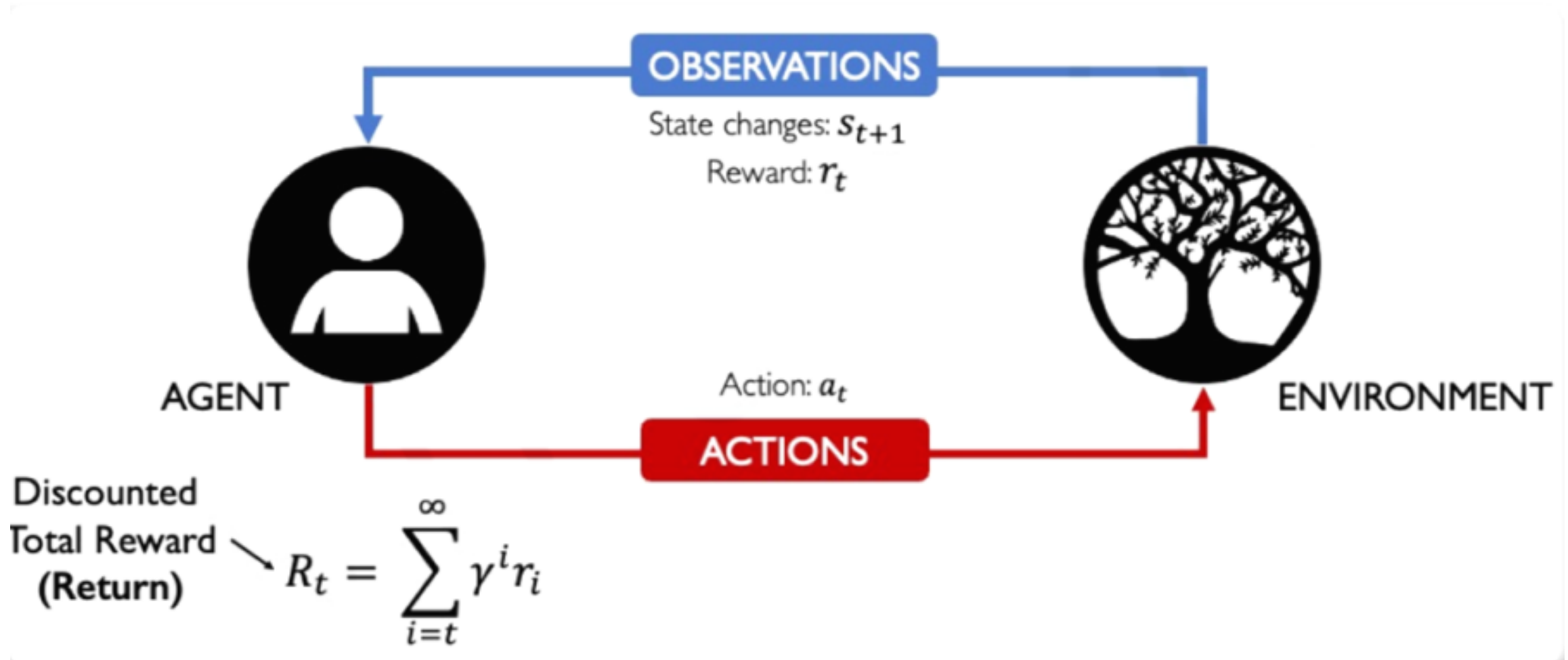
Reinforcement Learning: Concepts



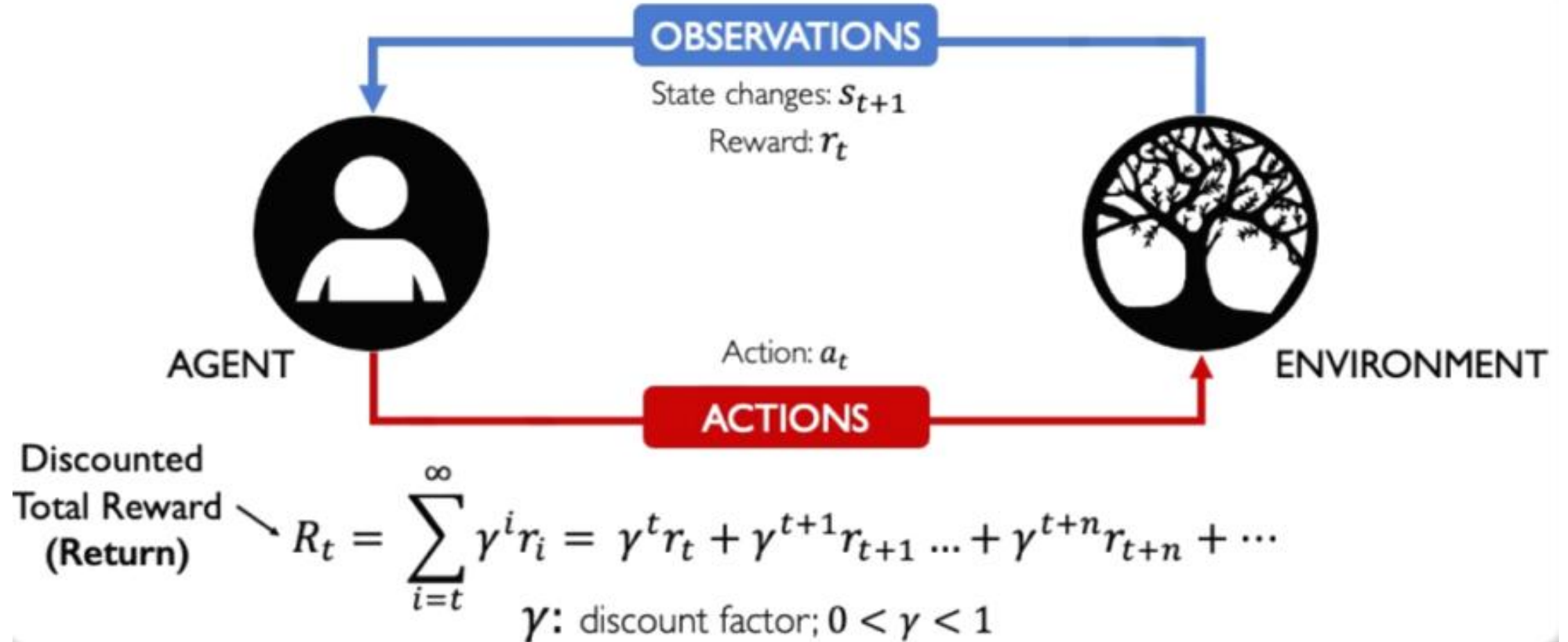
Reinforcement Learning: Concepts



Reinforcement Learning: Concepts



Reinforcement Learning: Concepts



Deep Reinforcement Learning Algorithms

Value Learning

Find $Q(s, a)$

$$a = \underset{a}{\operatorname{argmax}} Q(s, a)$$

Policy Learning

Find $\pi(s)$

Sample $a \sim \pi(s)$

Deep Reinforcement Learning Algorithms

Value Learning

Find $Q(s, a)$

$$a = \underset{a}{\operatorname{argmax}} Q(s, a)$$

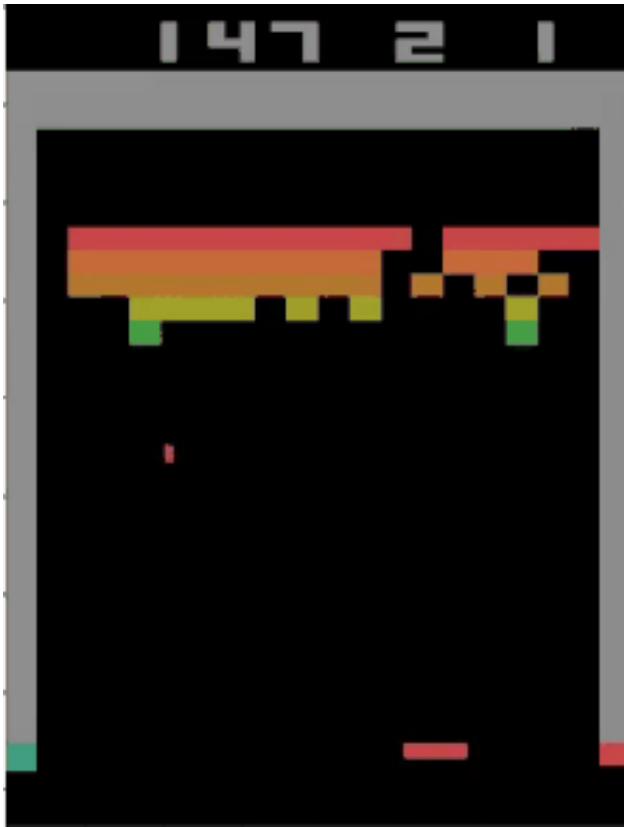
Policy Learning

Find $\pi(s)$

Sample $a \sim \pi(s)$

Deep Reinforcement Learning Algorithms

Atari Breakout



$$Q(s_t, a_t) = \mathbb{E}[R_t | s_t, a_t]$$

↑
↑
(state, action)

Strategy: the policy should choose an action that maximizes future reward

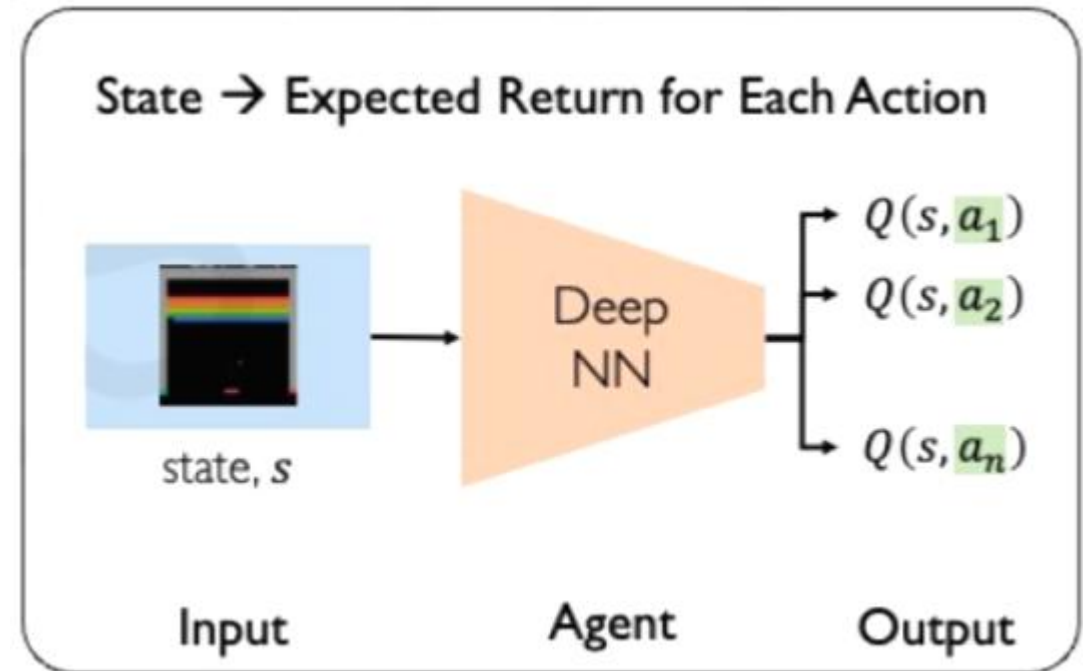
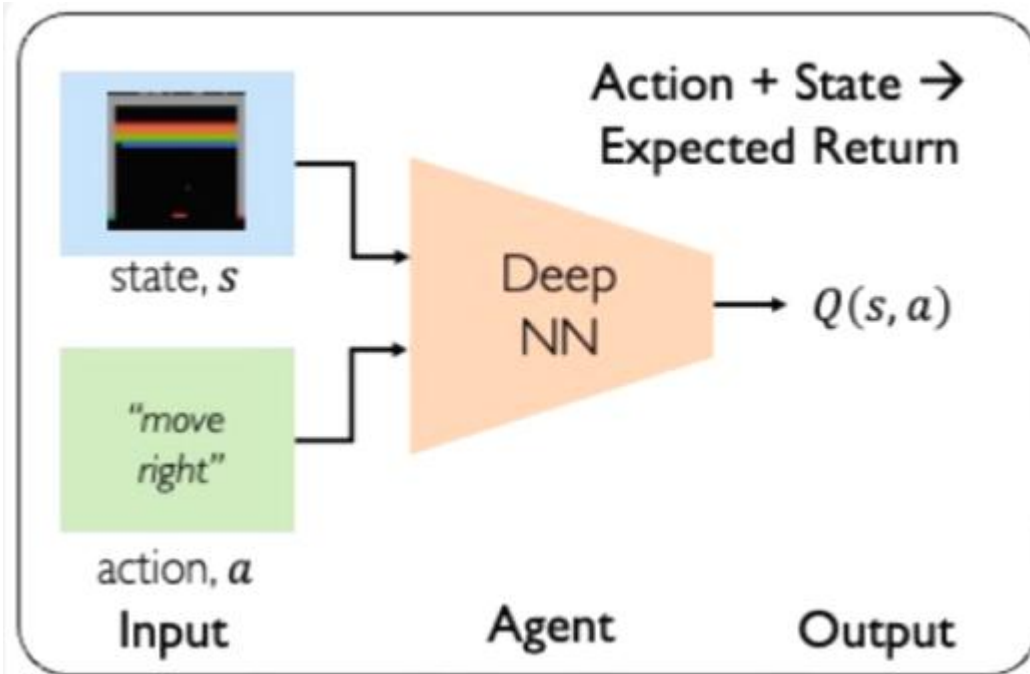
$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$$

Solution: Non Linear Function approximation

$$Q_w(s, a) \approx Q_\pi(s, a)$$

Deep Q Networks (DQN): Mnih et al. (2013/2015)

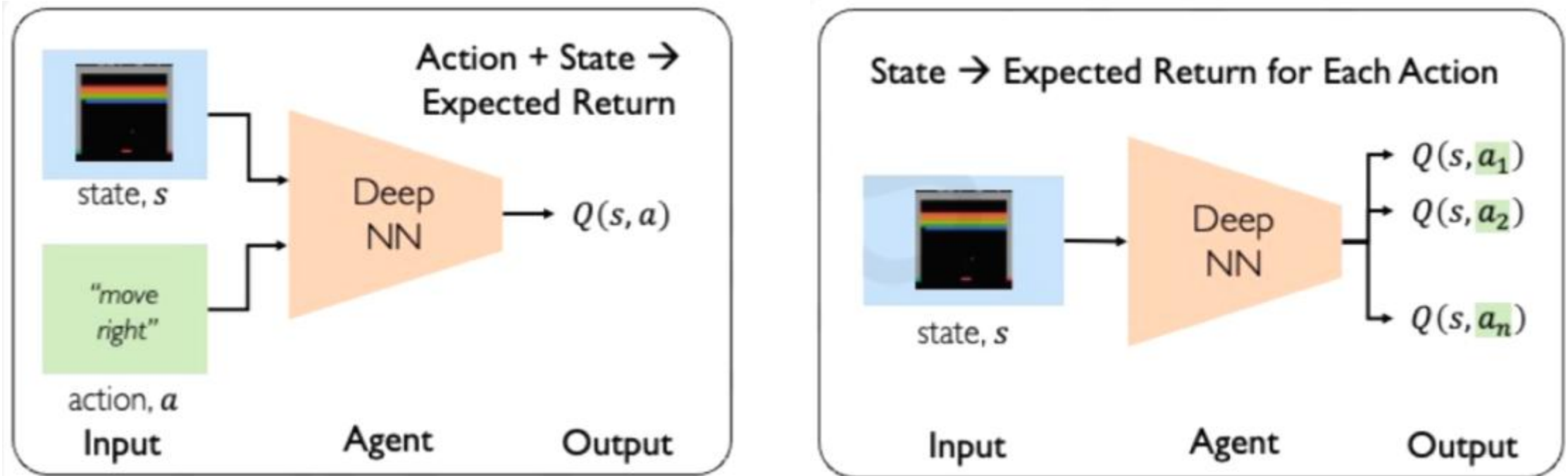
How can we use deep neural networks to model Q-functions?



What happens if we take all the best actions?
 Maximize target return \rightarrow train the agent

Deep Q Networks (DQN): Training

How can we use deep neural networks to model Q-functions?

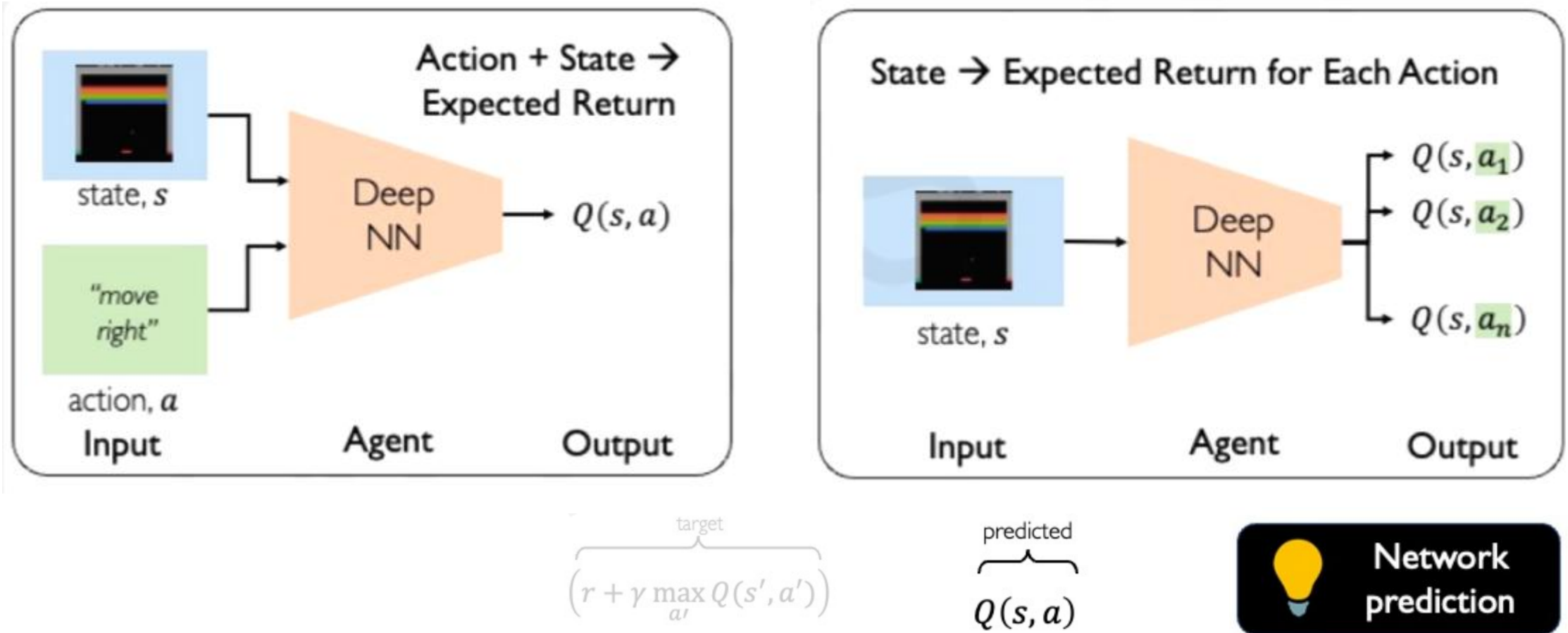


$$\overbrace{\left(r + \gamma \max_{a'} Q(s', a') \right)}^{\text{target}}$$

💡 Take all the best actions \rightarrow target return

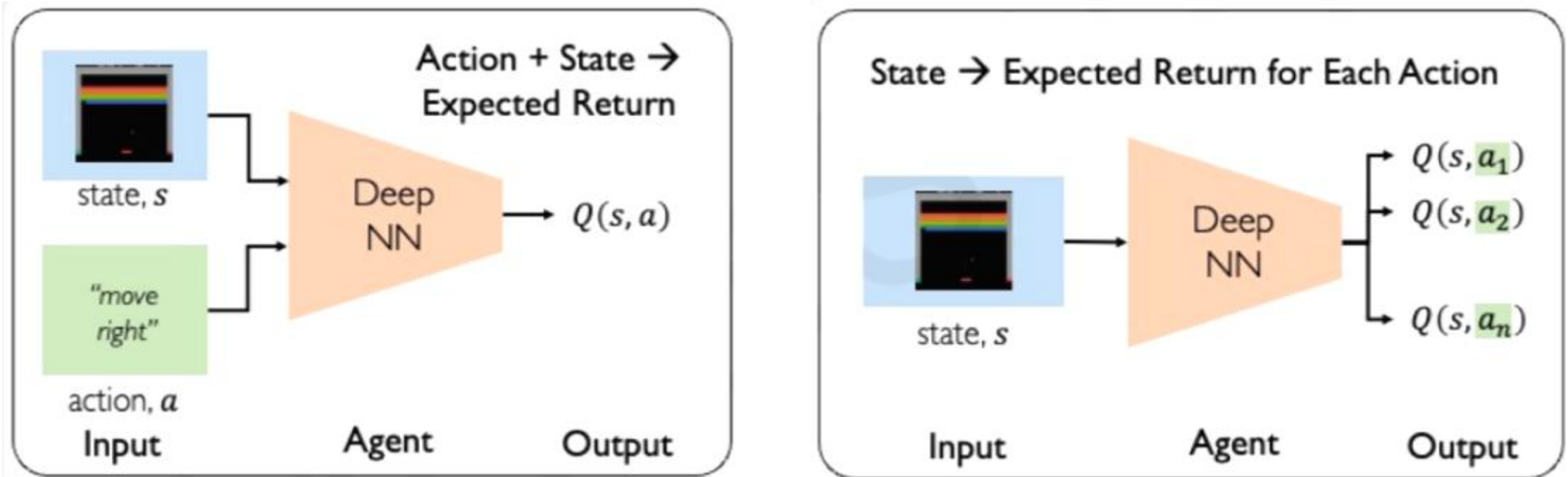
Deep Q Networks (DQN): Training

How can we use deep neural networks to model Q-functions?



Deep Q Networks (DQN): Training

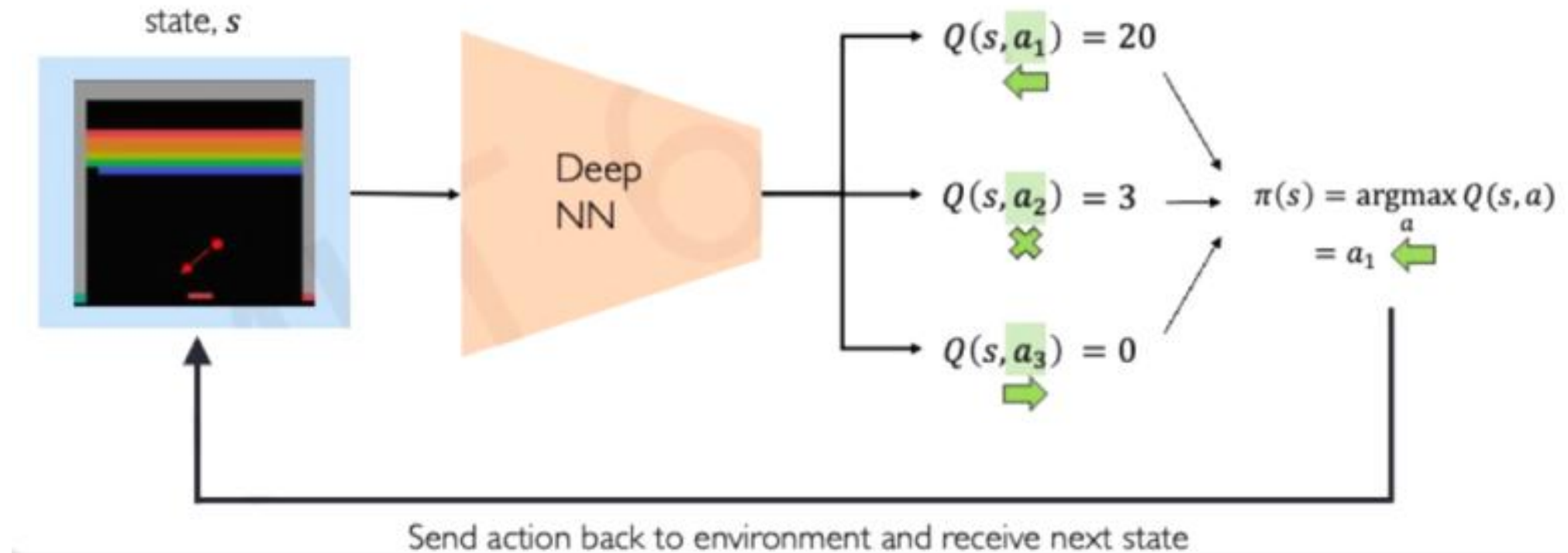
How can we use deep neural networks to model Q-functions?



$$\mathcal{L} = \mathbb{E} \left[\left\| \underbrace{\left(r + \gamma \max_{a'} Q(s', a') \right)}_{\text{target}} - \underbrace{Q(s, a)}_{\text{predicted}} \right\|^2 \right] \quad \text{Q-Loss}$$

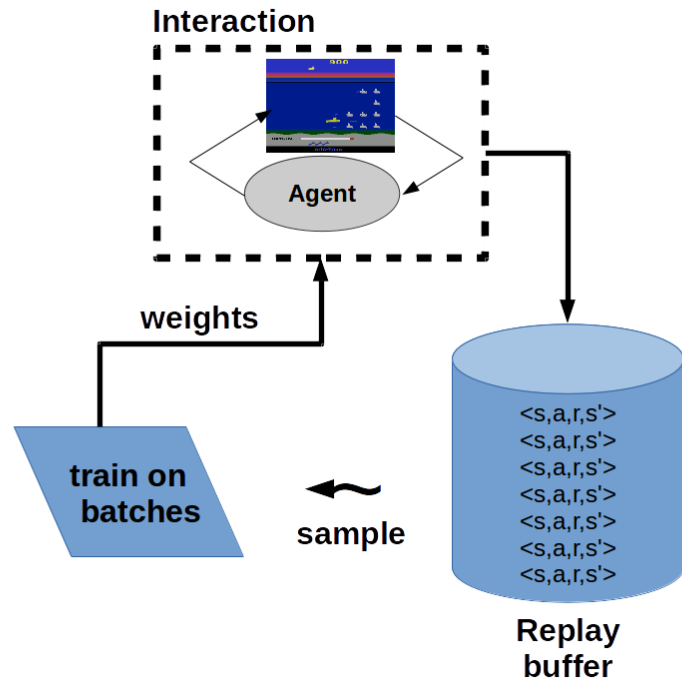
Deep Q Network Summary

Use NN to learn Q-function and then use to infer the optimal policy, $\pi(s)$



Deep Q Network with Experience Replay

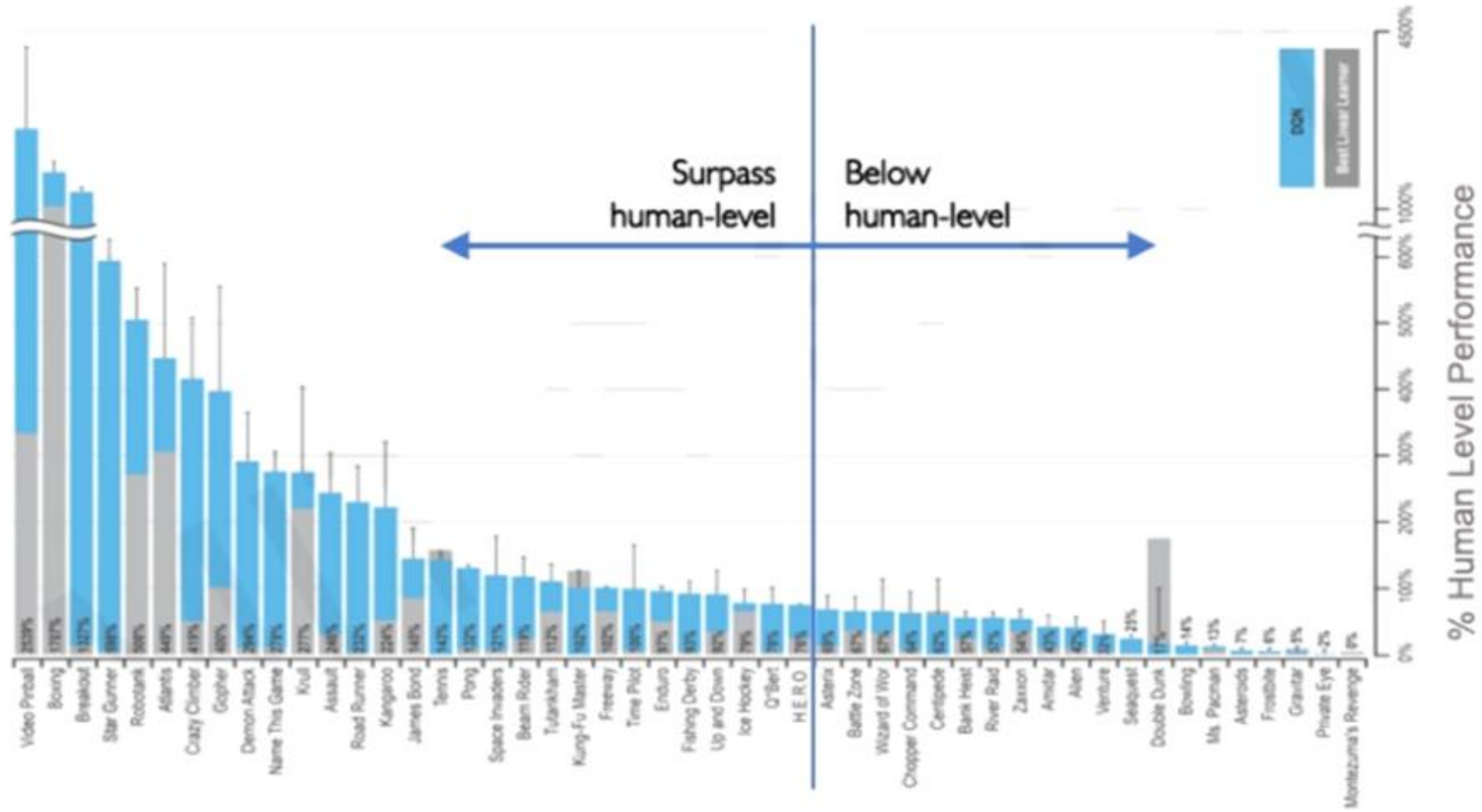
There's a powerful technique that you can use to improve sample efficiency for off-policy algorithms: Experience replay



Training with experience replay

- Play game, sample $\langle s, a, r, s' \rangle$.
- Update q-values based on $\langle s, a, r, s' \rangle$.
- Store $\langle s, a, r, s' \rangle$ transition in a buffer.
- If buffer is full, delete earliest data.
- Sample K such transitions from that buffer and update q-values based on them.

DQN Atari Results



Downsides of Q-learning

Complexity:

- Can model scenarios where the action space is discrete and small
- Cannot handle continuous action spaces

Flexibility:

- Policy is deterministically computed from the Q function by maximizing the reward → cannot learn stochastic policies

To address these, consider a new class of RL training algorithms: Policy gradient methods

Deep Reinforcement Learning Algorithms

Value Learning

Find $Q(s, a)$

$$a = \underset{a}{\operatorname{argmax}} Q(s, a)$$

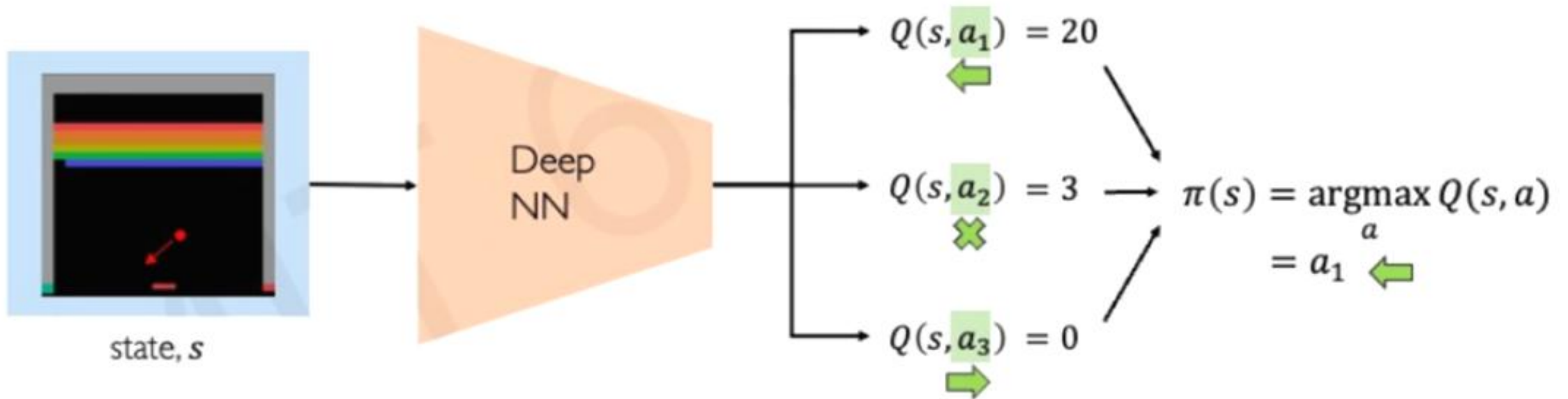
Policy Learning

Find $\pi(s)$

Sample $a \sim \pi(s)$

Deep Q Networks (DQN)

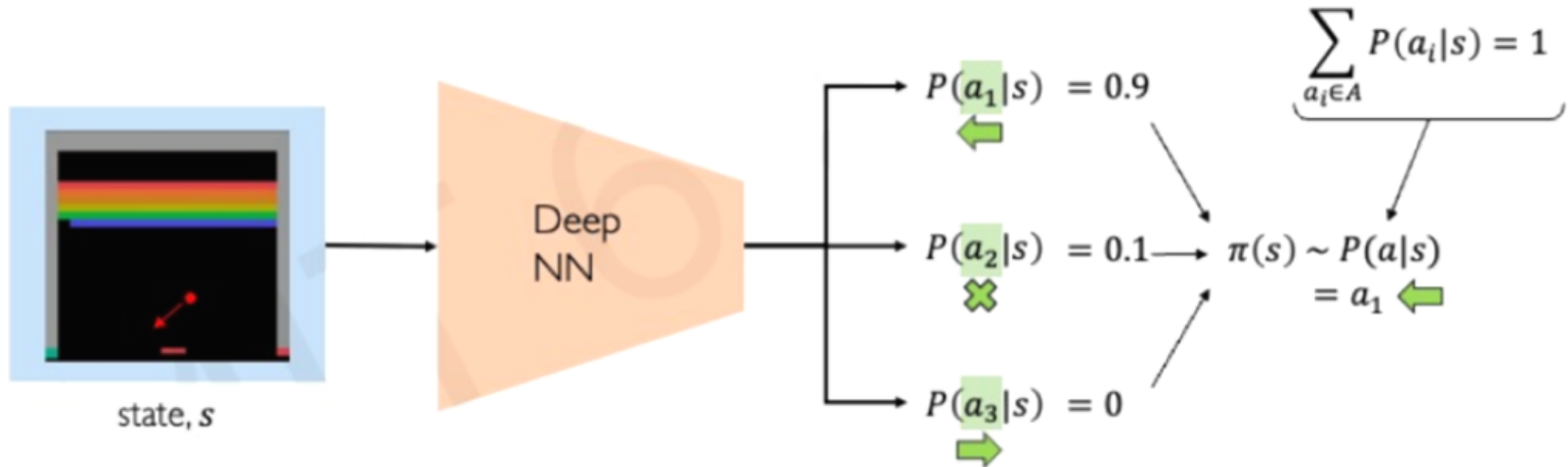
DQN: Approximate Q-function and use to infer the optimal policy, $\pi(s)$



Policy Gradient (PG): Key Idea

DQN: Approximate Q-function and use to infer the optimal policy, $\pi(s)$

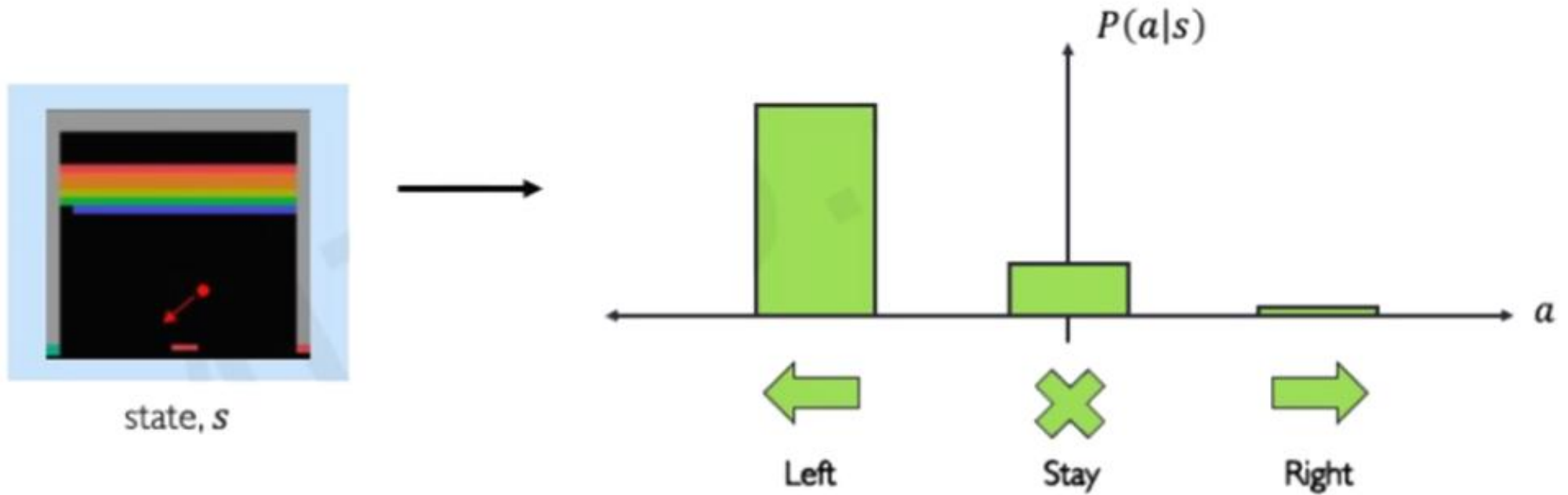
Policy Gradient: Directly optimize the policy $\pi(s)$



What are some advantages of this formulation?

Discrete vs Continuous Action Spaces

Discrete action space: which direction should I move? ← × →

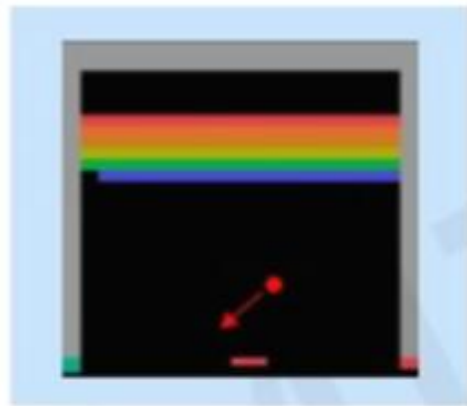


Discrete vs Continuous Action Spaces

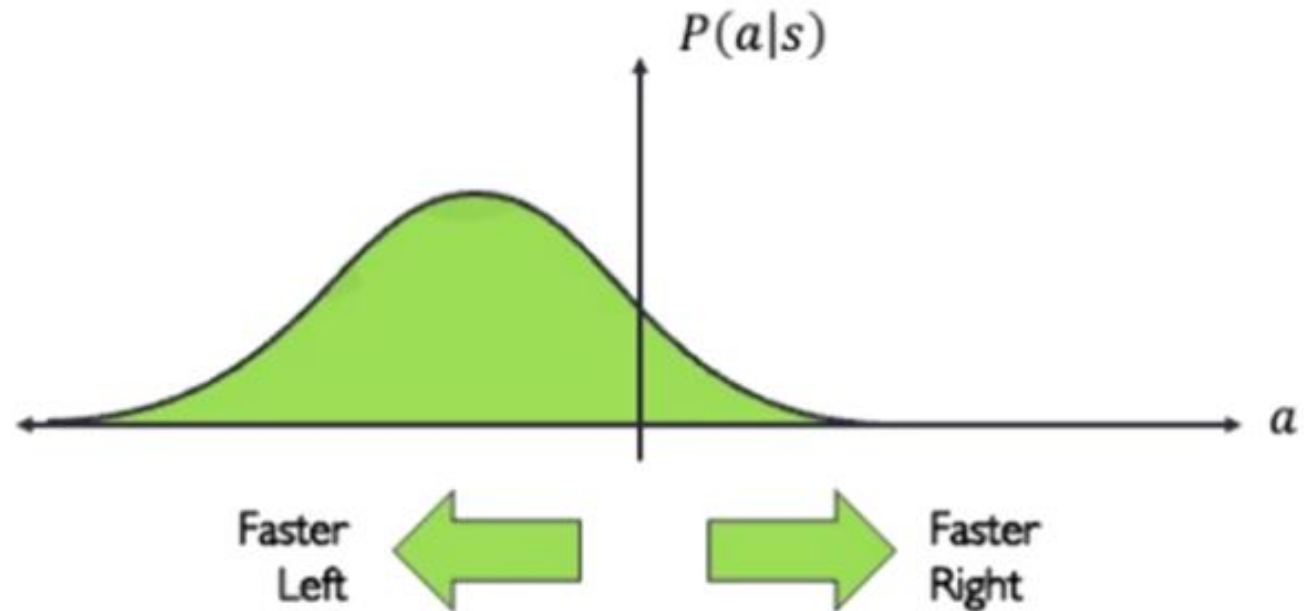
Discrete action space: which direction should I move?



Continuous action space: how fast should I move?

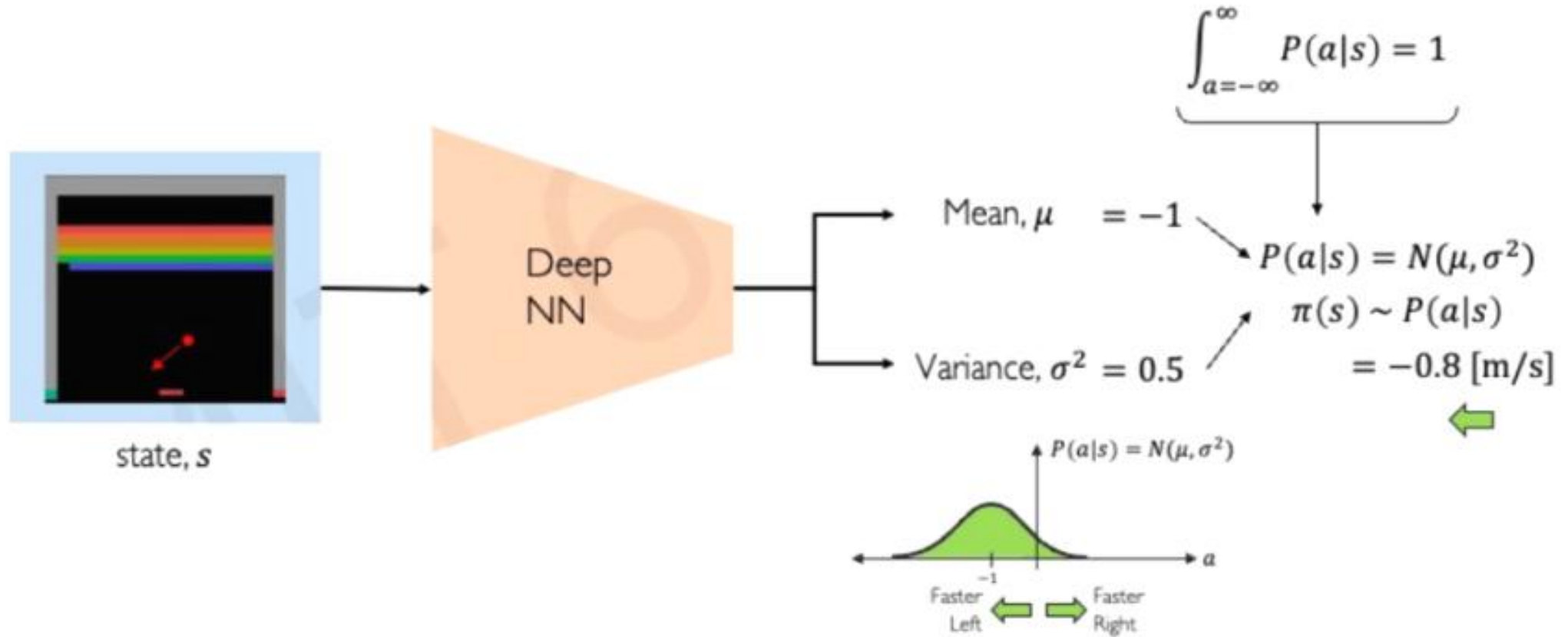


state, s



Policy Gradient (PG): Key Idea

Policy Gradient: Enables modeling of continuous action space



Training Policy Gradients: Case Study

Reinforcement Learning Loop:



Case Study – Self-Driving Cars

Agent: vehicle

State: camera, lidar, etc

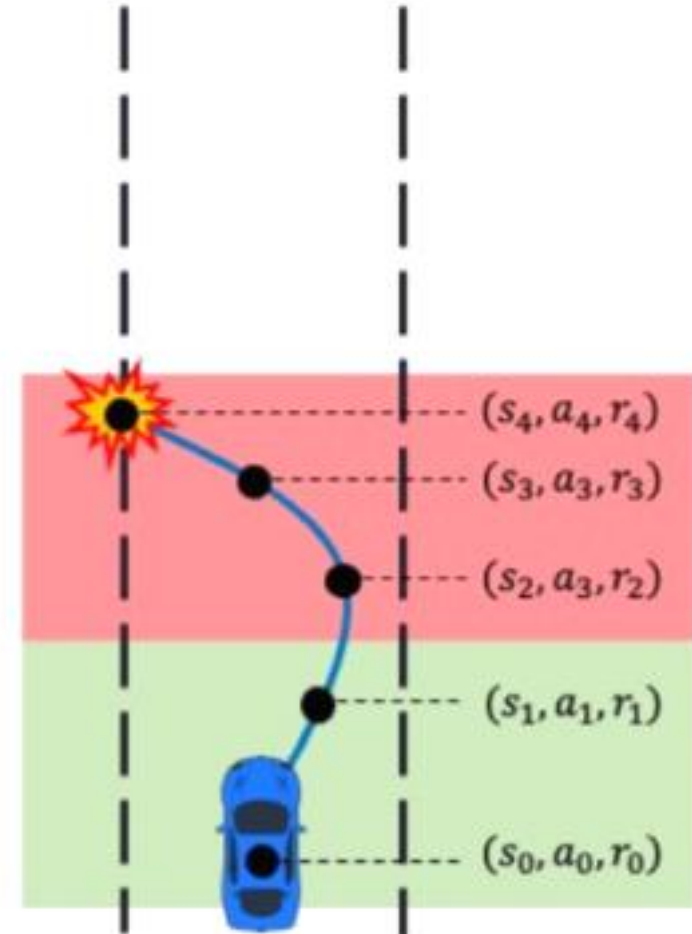
Action: steering wheel angle

Reward: distance traveled

Training Policy Gradients

Training Algorithm

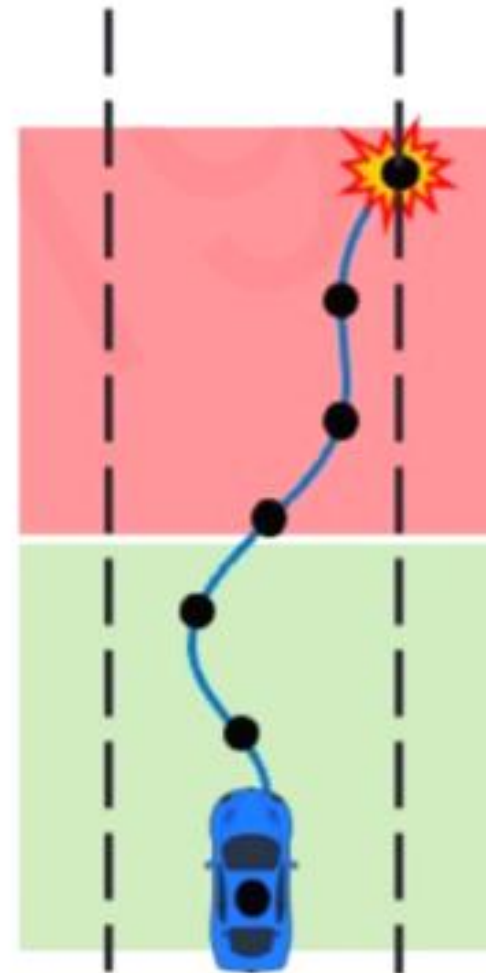
1. Initialize the agent
2. Run a policy until termination
3. Record all states, actions, rewards
4. Decrease probability of actions that resulted in low reward
5. Increase probability of actions that resulted in high reward



Training Policy Gradients

Training Algorithm

1. Initialize the agent
2. Run a policy until termination
3. Record all states, actions, rewards
4. Decrease probability of actions that resulted in low reward
5. Increase probability of actions that resulted in high reward



Training Policy Gradients

Training Algorithm

1. Initialize the agent
2. Run a policy until termination
3. Record all states, actions, rewards
4. Decrease probability of actions that resulted in low reward
5. Increase probability of actions that resulted in high reward



Training Policy Gradients

Training Algorithm

1. Initialize the agent
2. Run a policy until termination
3. Record all states, actions, rewards
4. Decrease probability of actions that resulted in low reward
5. Increase probability of actions that resulted in high reward



Training Policy Gradients

Training Algorithm

1. Initialize the agent
2. Run a policy until termination
3. Record all states, actions, rewards
4. Decrease probability of actions that resulted in low reward
5. Increase probability of actions that resulted in high reward

log-likelihood of action

$$\mathbf{loss} = -\log P(a_t | s_t) R_t$$

reward

Gradient descent update:

$$w' = w - \nabla \mathbf{loss}$$

$$w' = w + \nabla \log P(a_t | s_t) R_t$$

Policy gradient!

Policy Gradient Theorem

- ▶ The policy gradient approach also applies to (multi-step) MDPs
- ▶ Replaces reward R with long-term return G_t or value $q_\pi(s, a)$
- ▶ There are actually two policy gradient theorems (Sutton et al., 2000):
average return per episode & **average reward per step**

Policy Gradient Theorem(Episodic)

Theorem

For any differentiable policy $\pi_{\theta}(s, a)$, let d_0 be the starting distribution over states in which we begin an episode. Then, the policy gradient of $J(\theta) = \mathbb{E}[G_0 \mid S_0 \sim d_0]$ is

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T \gamma^t q_{\pi_{\theta}}(S_t, A_t) \nabla_{\theta} \log \pi_{\theta}(A_t | S_t) \mid S_0 \sim d_0 \right]$$

where

$$\begin{aligned} q_{\pi}(s, a) &= \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a] \end{aligned}$$

Policy Gradient Theorem(Average Reward)

Theorem

For any differentiable policy $\pi_{\theta}(s, a)$, the policy gradient of $J(\theta) = \mathbb{E}[R | \pi]$ is

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi} [q_{\pi_{\theta}}(S_t, A_t) \nabla_{\theta} \log \pi_{\theta}(A_t | S_t)]$$

Monte Carlo Policy Gradient (REINFORCE)

Algorithm 1 Episodic REINFORCE

- 1: Initialize policy network π_θ
 - 2: **for** $iteration = 1, 2, \dots, num_episodes$ **do**
 - 3: Generate an episode $S_0, A_0, R_1, S_1, \dots, S_{T-1}, A_{T-1}, R_t$ following π_θ
 - 4: **for** $t = 0, 1, 2, \dots, T - 1$ **do**
 - 5: $G_t = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$
 - 6: **end for**
 - 7: $L(\theta) = -\frac{1}{T} \sum_{t=0}^{T-1} G_t \ln \pi_\theta(A_t|S_t)$
 - 8: Update π_θ using $Adam(\nabla_\theta L(\theta))$
 - 9: **end for**
-

Monte Carlo Policy Gradient (REINFORCE) with baseline

Algorithm 2 Episodic REINFORCE with baseline

Initialize policy network π_θ

2: Initialize baseline network b_w

for $iteration = 1, 2, \dots, num_episodes$ **do**

4: Generate an episode $S_0, A_0, R_1, S_1, \dots, S_{T-1}, A_{T-1}, R_t$ following π_θ

for $t = 0, 1, 2, \dots, T - 1$ **do**

6: $G_t = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$

end for

8: $L(\theta) = -\frac{1}{T} \sum_{t=0}^{T-1} (G_t - b_w(S_t)) \ln \pi_\theta(A_t|S_t)$
 $L(w) = -\frac{1}{T} \sum_{t=0}^{T-1} (G_t - b_w(S_t))^2$

10: Update π_θ using $Adam(\nabla_\theta L(\theta))$
 Update b_w using $Adam(\nabla_w L(w))$

12: **end for**

Deep RL: Value and Policy Based Summary

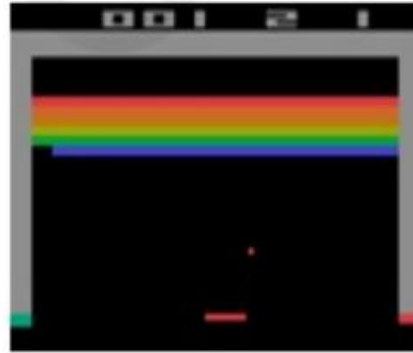
Foundations

- Agents acting in environment
- State-action pairs \rightarrow maximize future rewards
- Discounting



Q-Learning

- Q function: expected total reward given s, a
- Policy determined by selecting action that maximizes Q function



Policy Gradients

- Learn and optimize the policy directly
- Applicable to continuous action spaces

