# REINFORCEMENT LEARNING

**Hyundai Motors Boot-campus**

UKJO HWANG, DOHYEOK KWON

**INFORMATION AND INTELLIGENCE SYSTEMS LAB.**

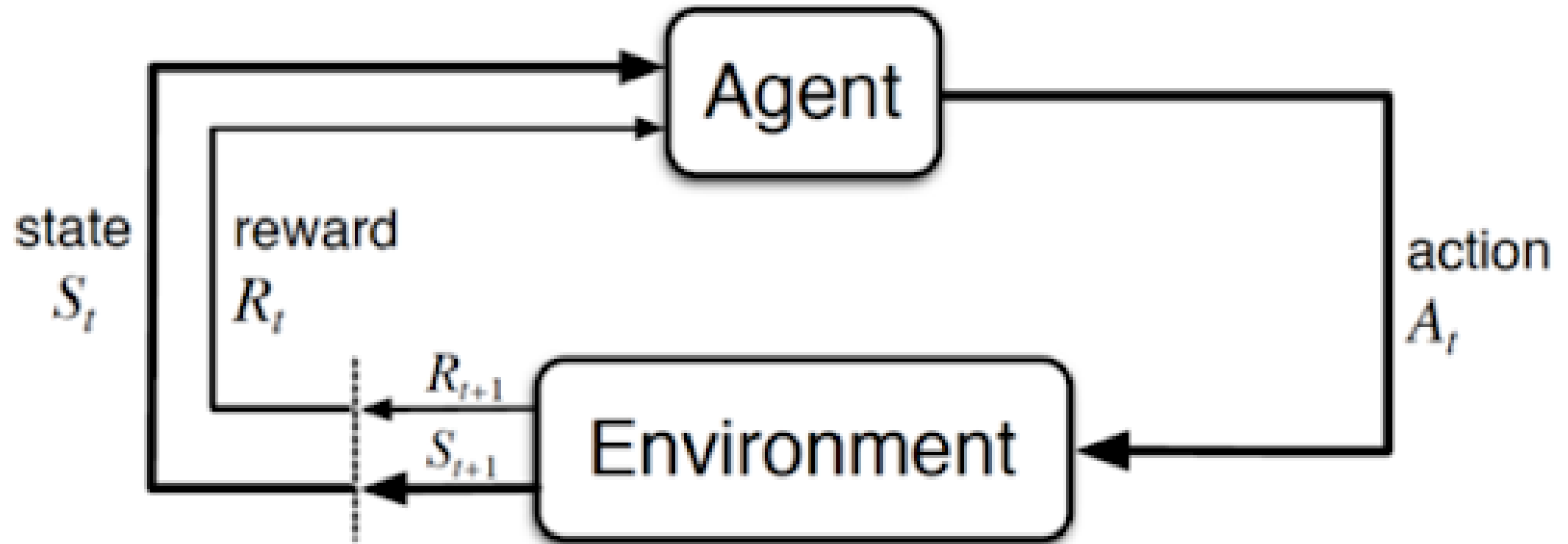ELECTRONIC ENGINEERING, HANYANG UNIVERSITY

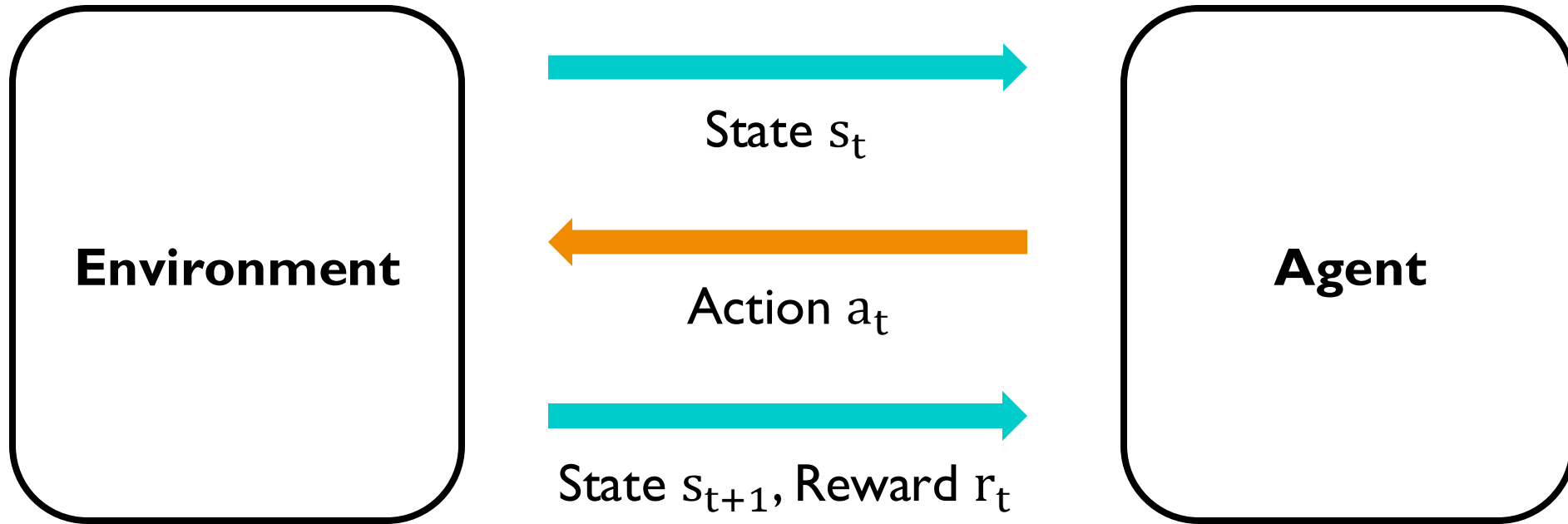**July 7, 2025**

# GOAL

- **How to implement deep RL algorithms**

    - How to utilize environments using the OpenAI Gym API

    - How to train neural networks for reinforcement learning tasks

    - How to train agents using DQN and DDPG algorithms
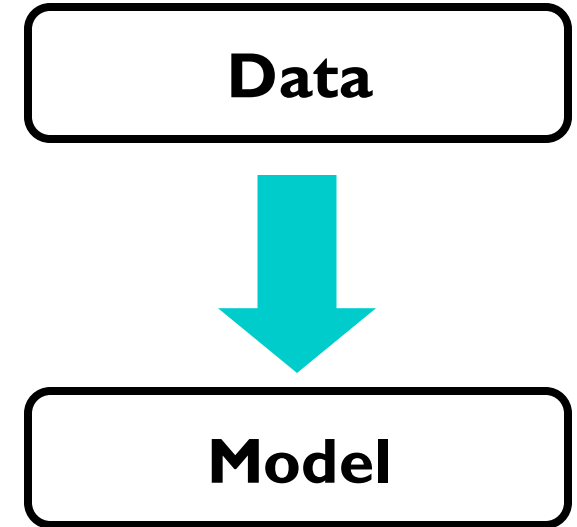
# ENVIRONMENT AGENT INTERACTION



- The environment provides the agent with information about the current **state** and **reward**.

- The agent selects an **action** based on its policy.

- The environment transitions to a new state according to its **state transition probabilities.**
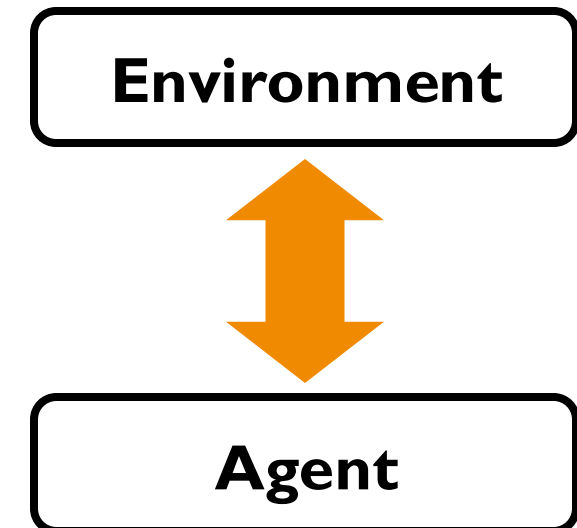
# SUPERVISED LEARNING VS REINFORCEMENT LEARNING

- **Supervised Learning**

  – The data is provided in advance.

  – The distribution and quality of data are given beforehand.

```
┌──────────────┐
│     Data     │
└──────────────┘
       ↓
┌──────────────┐
│    Model     │
└──────────────┘
```

- **Reinforcement Learning**
  – The data is **generated through interactions** between the

    agent and the environment during the learning process

  – The distribution and quality of data continuously change according

    to the action (policy).

```
┌──────────────┐
│ Environment  │
└──────────────┘
       ↕
┌──────────────┐
│    Agent     │
└──────────────┘
```

# HOW TO IMPLEMENT RL ALGORITHMS

**Environment**

– What is the shape and structure of the data?

– What functions and APIs does the environment provide?

**Agent**

– How to design and implement the model (e.g., Neural Network)

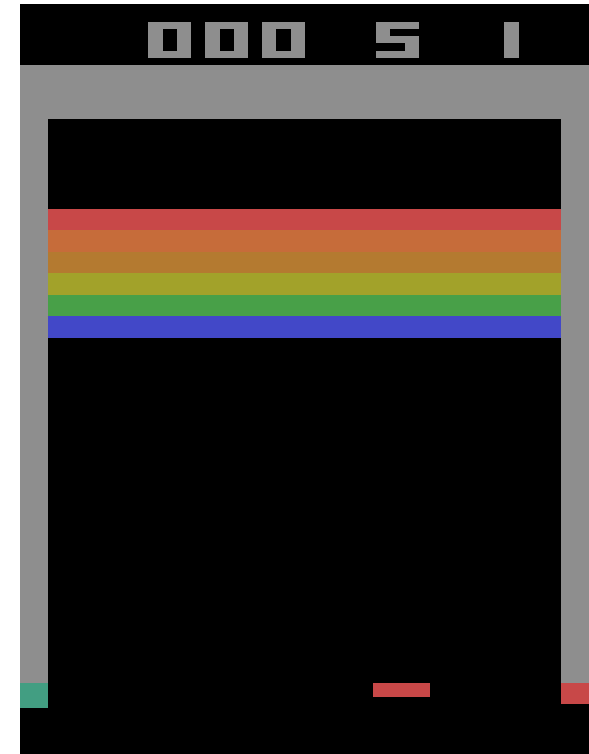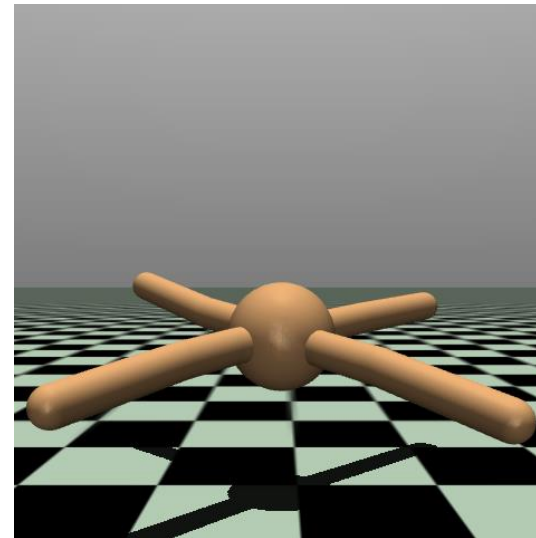– How to train the model using the provided data

**Combination**

– How the agent interacts with the environment

– How to implement and apply specific RL algorithms

# OpenAI Gymnasium

- **OpenAI Gymnasium**

  - Provides a wide range of environments for RL.

  - It's important to understand the basic interface of Gymnasium.

  - To implement a new environment, it's best to follow the Gymnasium API standards.

  - https://gymnasium.farama.org/

# SETUP

- ## I. Create Anaconda virtual environment

  - conda create –n hyundai_rl python=3.10

  - conda activate hyundai_rl

- ## II. Install packages

  - pip install gymnasium matplotlib ipykernel torch

# I. Q-Learning

- **Q-Value Function(or State-Action Value Function)** $Q(S, A)$
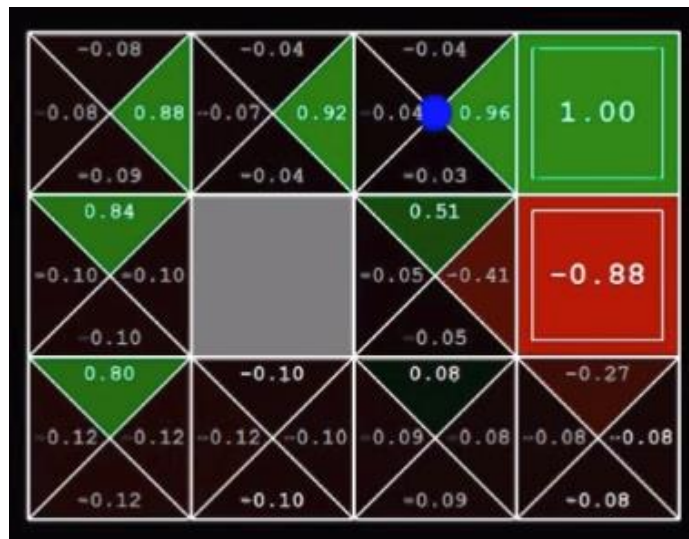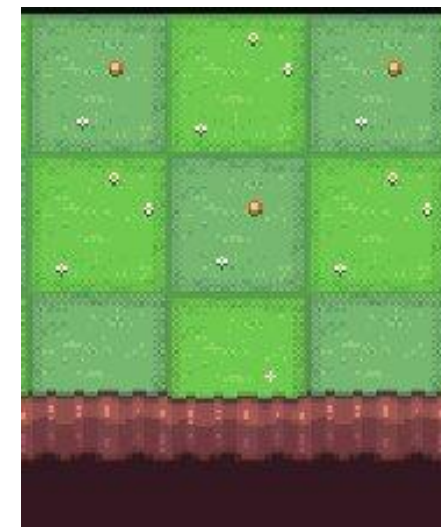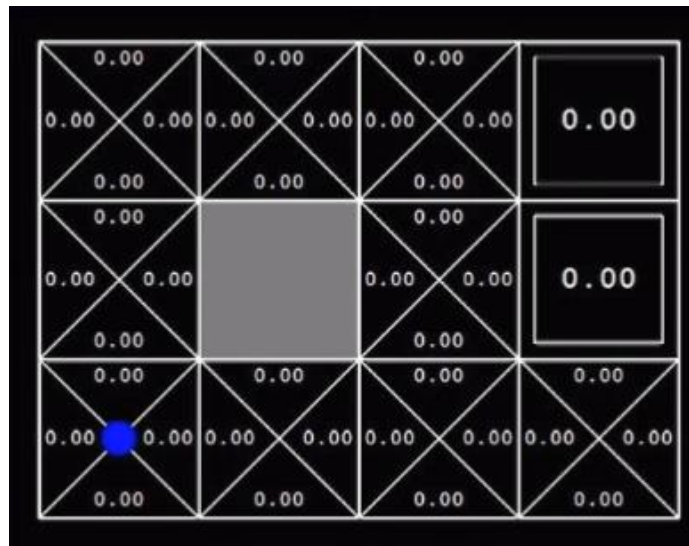  - 특정상태 $S(State)$에서 특정 행동 $A(Action)$를 선택했을 때 기대되는 누적 보상.

$$Q(S, A) = \mathbb{E}[R_t | s_t = S, a_t = A] = \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t = S, a_t = A\right]$$

  - 학습을 통해서 모든 state와 action에 대한 Q 값을 찾으면, 특정 state에서 어떤 action을 선택하는 것이 가장 큰 보상을 받을 수 있을지 아는 것과 같음.

# BACKGROUND

- **Q-Table**

    – Q-Value를 저장하는 간단한 방법.

    – 모든 state, action pair에 대해서 Q-value를
      저장해야 함.

    – Dimension: $S \times A$





|       | Left | Right | Up | Down |
|-------|------|-------|----|------|
| (1,1) | 0    | 0     | 0  | 0    |
| (1,2) | 0    | 0     | 0  | 0    |
| (1,3) | 0    | 0     | 0  | 0    |
| (2,1) | 0    | 0     | 0  | 0    |
| (2,2) | 0    | 0     | 0  | 0    |
| (2,3) | 0    | 0     | 0  | 0    |
| (3,1) | 0    | 0     | 0  | 0    |
| (3,2) | 0    | 0     | 0  | 0    |
| (3,3) | 0    | 0     | 0  | 0    |

# Q-LEARNING

**Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:

    Initialize $S$

    Loop for each step of episode:

        Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)

        Take action $A$, observe $R$, $S'$

        $Q(S, A) \leftarrow Q(S, A) + \alpha \left[ R + \gamma \max_a Q(S', a) - Q(S, A) \right]$

        $S \leftarrow S'$

    until $S$ is terminal

# Q-LEARNING

- **1. Q-Table 초기화**

$$Q(S, A) = 0 \quad \forall S, A$$

- **2. $\epsilon$-greedy policy**
  - 탐색(Exploration)과 활용(Exploitation)을 적절히 조합하여 다음 action을 선택해야 함.
  - $\epsilon$의 확률로 random action을 선택하고, $1 - \epsilon$의 확률로 현재 state에서 가장 높은 Q-value를 갖는 action 선택
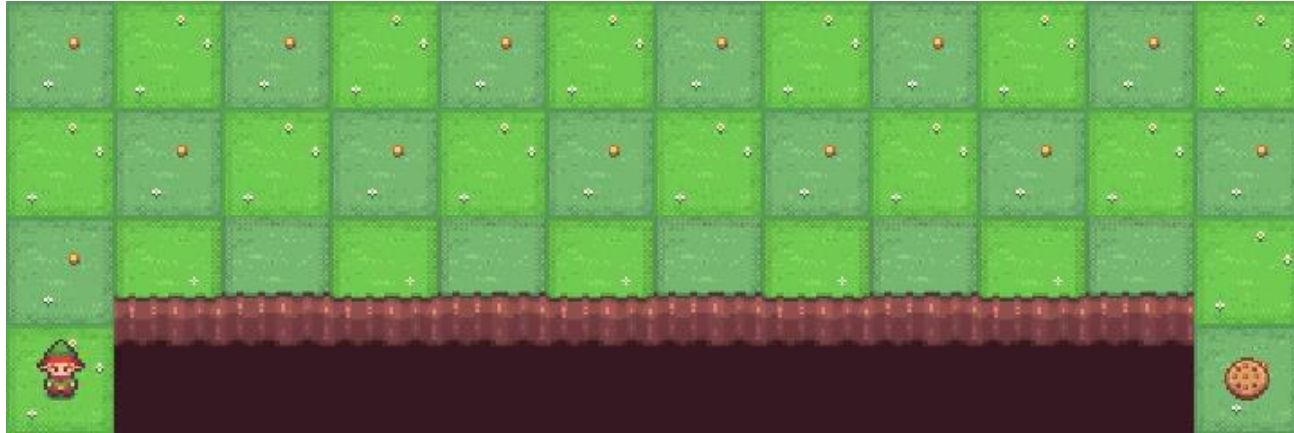
- **3. Q-value update**
  - Bellman equation을 따라 Q-value를 update.

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_{a'} Q(S', a') - Q(S, A)]$$

  - $\alpha$: Learning rate / $\gamma$ : Discount factor $\Rightarrow$ 모두 [0, 1] 사이의 실수 값
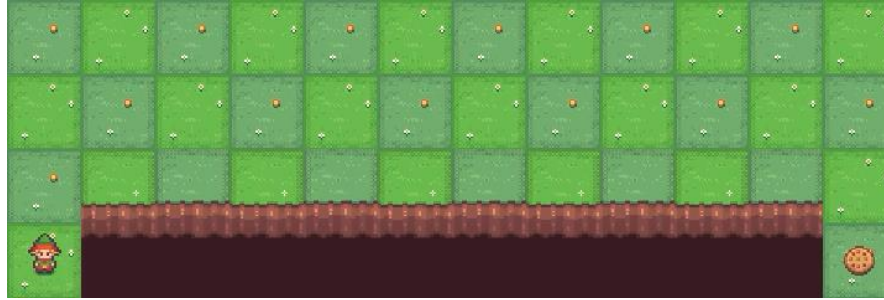
# ENV: CLIFF WALKING



- **Action Space:** Discrete(4) | **Observation Space:** Discrete(48)

  - **[Action]** 0: Move up | 1: Move right | 2: Move down | 3: Move left

- **Starting State:** [36] (3, 0) | **Episode End:** [47] (3, 11)

- **Reward**

  - Each time step incurs -1 reward, unless the player stepped into the cliff, which incurs -100 reward.
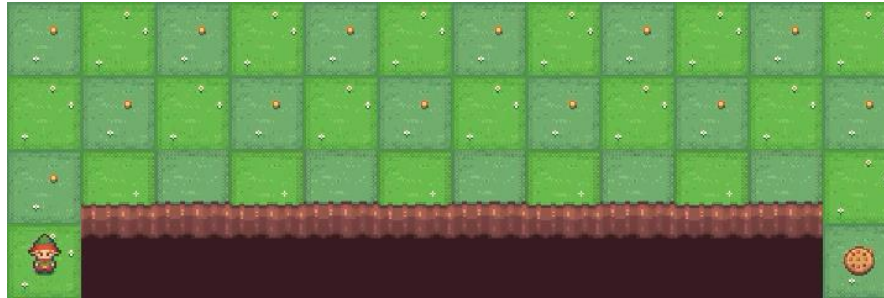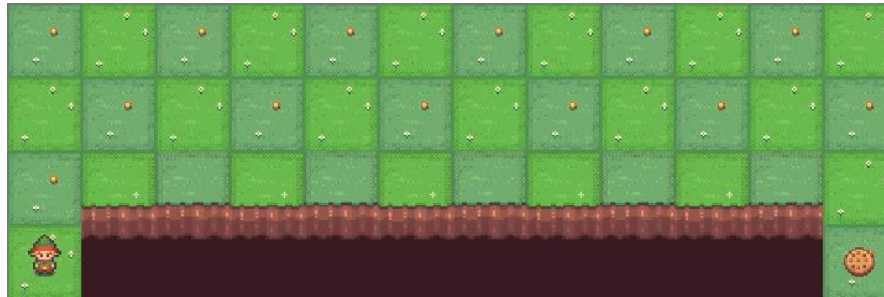
# ENV: CLIFF WALKING

- **[Training]**

  – Episode 0~10 (Average steps: 240.0)



  – Episode 100~110 (Average steps: 38.5)



  – Episode 200~210 (Average steps: 17.1)

# EXPERIMENTAL EXTENSIONS

- **Exploration probability** 가 더 높은 값이라면?

  – Increase 'exploration rate': 0.1→ **0.4**

  – Decrease 'number of episodes': 300 → **150**


- 유사한 환경인 `FrozenLake`에서 학습에 **더 많은 step**이 필요한 이유?

  – env = gym.make('FrozenLake-v1', render_mode='rgb_array', is_slippery=False)

  – Increase episodes from 150 to 2000


- **Stochastic transition**의 영향 확인

  – Set 'is_slippery = True'

# EXPERIMENTAL EXTENSIONS

- **Exploration probability** 가 더 높은 값이라면?

- 유사한 환경인 `FrozenLake`에서 학습에 **더 많은 step**이 필요한 이유?

- **Stochastic transition**의 영향 확인



1500

2000

Slippery

# II. DQN

# DEEP NEURAL NETWORKS

- If the state and action spaces are large, it becomes impractical to implement the Q-function using a table.

- Therefore, deep neural networks are widely used to approximate the Q-function.

# Q-FUNCTION

- The output of the DNN is the Q-value corresponding to a given state-action pair.

- Therefore, it is natural to design the DNN to take the state and action as input and output the corresponding Q-value.

State s →

Action a →

Q Network

$Q(s, a)$ →

# Q-Function

- However, in practice, the DNN is often designed to take only the state as input and output the Q-values for every possible action.

- The feasibility of this implementation relies on the discreteness of the action space.

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right) \nabla_{\theta_i} Q(s, a; \theta_i) \right]$$

$Q(s, a_1)$

$Q(s, a_2)$

State s

Q Network

$Q(s, a_3)$

# IMPLEMENTATION DETAILS

- Since DQN is an off-policy algorithm, it utilizes a replay buffer to store samples generated while interacting with the environment.

- This design choice is crucial because it helps to reduce strong correlations between samples.



Store experience tuples

$(s_t^{(i)}, a_t^{(i)}, r_{t+1}^{(i)}, s_{t+1}^{(i)})$

$(s_t^{(1)}, a_t^{(1)}, r_{t+1}^{(1)}, s_{t+1}^{(1)})$
$(s_t^{(2)}, a_t^{(2)}, r_{t+1}^{(2)}, s_{t+1}^{(2)})$
$(s_t^{(3)}, a_t^{(3)}, r_{t+1}^{(3)}, s_{t+1}^{(3)})$

Sample minibatch (uniformly) for training

$\left\{(s_t^{(k)}, a_t^{(k)}, r_{t+1}^{(k)}, s_{t+1}^{(k)})\right\} \sim U(D)$

Replay Buffer (D)

# IMPLEMENTATION DETAILS

- DQN employs a target network that gradually updates to follow the parameters of the original network.

- This mechanism contributes to stabilizing the training.

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s,a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right) \nabla_{\theta_i} Q(s, a; \theta_i) \right].$$

# DQN Algorithm

---
**Algorithm 1** Deep Q-learning with Experience Replay
---

Initialize replay memory $\mathcal{D}$ to capacity $N$

Initialize action-value function $Q$ with random weights

**for** episode $= 1, M$ **do**

    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

    **for** $t = 1, T$ **do**

        With probability $\epsilon$ select a random action $a_t$

        otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$

        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$

$$\text{Set } y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$$

        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

    **end for**

**end for**

---

# III. Actor-Critic Algorithm

# BACKGROUND

- **Policy-based reinforcement learning**

    - 기존 Value-based algorithm이 어떤 state에서 각 action의 가치를 계산하여 행동했다면, Policy-based algorithm은 각 state에서의 행동 정책(확률)을 학습하는 방식.

        › DQN은 Network의 output으로 특정 (s, a)의 Q-Value를 출력

        › Policy Network는 각 현재 state에서 각 action을 수행할 확률 $\pi(a|s) = P(A_t = a \mid S_t = s)$를 출력.

- **Policy gradient**

    - Objective function: $J(\theta) = V_{\pi_\theta}(s_0)$ (To maximize it, $\theta_{t+1} = \theta_t + \alpha \nabla_\theta J(\theta)$)

        › $V_{\pi_\theta}(s) = \mathbb{E}_{\pi_\theta}\left[R_{t+1} + \gamma V_{\pi_\theta}(S_{t+1}) \mid S_t = s\right]$ (Expected return)

    - **Policy gradient theorem: $\boldsymbol{\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[r(\tau) \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t)]}$**

- **Policy gradient**

  - **Policy gradient theorem:** $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[r(\tau) \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t)]$
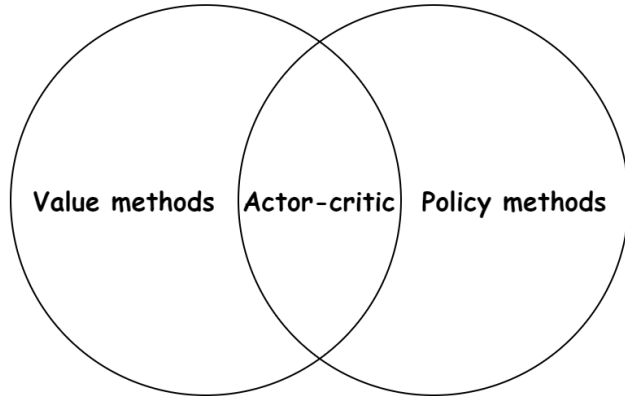
$$
\begin{aligned}
\nabla_\theta J(\theta) &= \mathbb{E}_{\pi_\theta}\left[\sum_{t=0}^{T-1} G_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)\right] && \text{REINFORCE (Monte Carlo PG)} \\
&= \mathbb{E}_{\pi_\theta}\left[\sum_{t=0}^{T-1} (G_t - V_\phi(s_t)) \nabla_\theta \log \pi_\theta(a_t \mid s_t)\right] && \text{REINFORCE with baseline} \\
&= \mathbb{E}_{\pi_\theta}\left[Q_\phi(s,a) \nabla_\theta \log \pi_\theta(a \mid s)\right] && \text{Q-value Actor-Critic} \\
&= \mathbb{E}_{\pi_\theta}\left[A_{\phi_1,\phi_2}(s,a) \nabla_\theta \log \pi_\theta(a \mid s)\right] && \text{Advantage Actor-Critic} \\
&= \mathbb{E}_{\pi_\theta}\left[(r + \gamma V_\phi(s') - V_\phi(s)) \nabla_\theta \log \pi_\theta(a \mid s)\right] && \text{TD Actor-Critic}
\end{aligned}
$$

Critic (value function)   Actor (policy)

# ACTOR-CRITIC



Value methods | Actor-critic | Policy methods

- **Actor (Policy Network)**
  - 각 state $s$에서 다음 액션에 대한 확률(정책- $\pi(s|a)$)을 근사하는 Network
  - Loss: $J(\theta)$

- **Critic (Value Network)**
  - 각 state $s$에서의 가치(Value function)를 근사하는 Network.
  - Loss: MSE between $V(s)$ and TD target $(r + \gamma V(s'))$

$01:$ **Input** : Initial Actor policy parameters $\pi(a|s, \theta)$
　　　　: Initial Critic V-function parameters $V(s, \psi)$
$02:$ **Parameters** : actor learning rate $\alpha^\theta > 0$, critic learning rate $\alpha^\psi > 0$
$03:$ trial_step $= T$
$04:$ **For** episode $= 1, M$ **do**
$05:$ 　　done $= False$
$06:$ 　　Reset environment state.
$07:$ 　　**For** $t = 1, T$ **do** ($or$ **While not** $done$)
$08:$ 　　　　Observe state $s$ and select action $a = \mu_\theta(s)$, 　　　# discrete model

$09:$ 　　　　Execute $a$ in the environment
$10:$ 　　　　Observe next state $s'$, reward $r$, and done signal $d$ to indicate whether $s'$ is terminal
$11:$ 　　　　Reset gradient $d\theta$ and $d\psi$ to 0
$12:$ 　　　　Calculate the TD

$$TD \leftarrow r + \gamma(1-d)V(s', \psi) - V(s, \psi) \quad \text{(if } s' \text{ is terminal, then } V(s', \psi) \doteq 0)$$

$$\text{where, } Q_{expected} \leftarrow r + \gamma(1-d)V(s', \psi)$$

$13:$ 　　　　Accumulate the policy gradient using the critic:
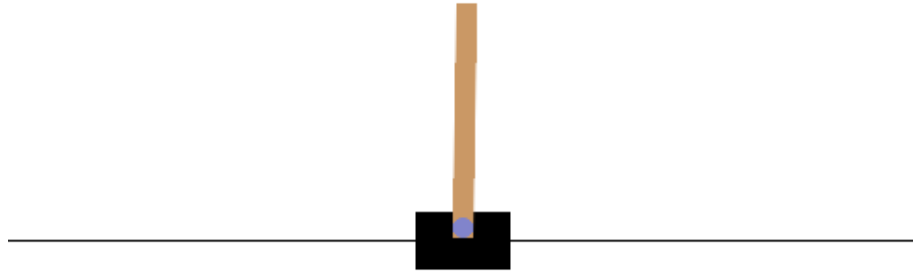$$d\theta \leftarrow d\theta + \nabla_\theta \log \pi_\theta(s_t, a_t)(TD)$$
$14:$ 　　　　Accumulate the critic gradient:
$$d\psi \leftarrow d\psi + \nabla_\psi (TD)^2$$

$15:$ 　　　　Update the actor and the critic with the accumulated gradients using gradient descent or similar :

$$\theta \leftarrow \theta + \alpha^\theta d\theta \qquad \psi \leftarrow \psi + \alpha^\psi d\psi$$

$16:$ 　　　　$s \leftarrow s'$
$17:$ 　　**End For**
$18:$ **End For**
$19:$

# ENV: CART POLE



| Num | Observation | Min | Max |
|-----|-------------|-----|-----|
| 0 | Cart Position | -4.8 | 4.8 |
| 1 | Cart Velocity | -Inf | Inf |
| 2 | Pole Angle | ~ -0.418 rad (-24°) | ~ 0.418 rad (24°) |
| 3 | Pole Angular Velocity | -Inf | Inf |

- **Action Space:** Discrete(2)  |  **Observation Space:** Continuous(4, )

  - **[Action]**  0: Push cart to the left  |  1: Push cart to the right

- **Episode End:** Angle $\pm12°$ / End of display / 500 steps

- **Reward**

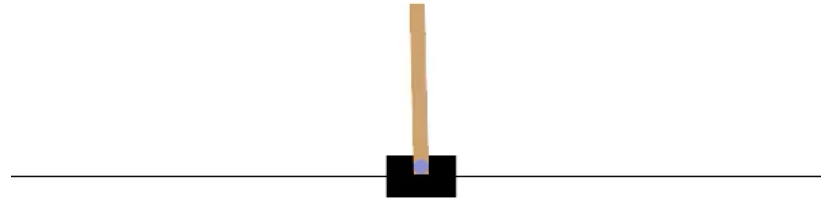  - A +1 reward is given for every step, including the terminal step.
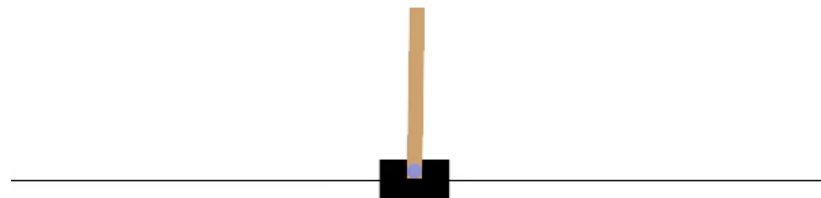
# ENV: CART POLE

- **[Training]**

  – Episode 100~200 (Average steps: 60.6)

  – Episode 200~210 (Average steps: 144.6)

  – Episode 500~510 (Average steps: 282.0)

# IV. DDPG

# DQN vs DDPG

- DQN

  - Designed for environments with discrete action spaces

  - A Q-network alone is sufficient for learning the policy

  $$\pi(s) = \max_{a \in \mathcal{A}} Q^{\pi}(s, a)$$

  $$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q^*(s', a') \Big| s, a \right]$$

- DDPG

  - Suitable for environments with continuous action spaces

  - Requires both a policy network (Actor) and a Q-network (Critic)

  $$Q^{\mu}(s_t, a_t) = \mathbb{E}_{r_t, s_{t+1} \sim E} \left[ r(s_t, a_t) + \gamma Q^{\mu}(s_{t+1}, \mu(s_{t+1})) \right]$$

# DETERMINISTIC POLICY

- Deterministic policies simplify the computation of gradients during the optimization of the policy network.

$$\nabla_{\theta^\mu} J \approx \mathbb{E}_{s_t \sim \rho^\beta} \left[ \nabla_{\theta^\mu} Q(s, a | \theta^Q) |_{s=s_t, a=\mu(s_t | \theta^\mu)} \right]$$
$$= \mathbb{E}_{s_t \sim \rho^\beta} \left[ \nabla_a Q(s, a | \theta^Q) |_{s=s_t, a=\mu(s_t)} \nabla_{\theta_\mu} \mu(s | \theta^\mu) |_{s=s_t} \right]$$

# CRITIC

- The critic network receives the state and action as input and produces the corresponding Q-value as output.

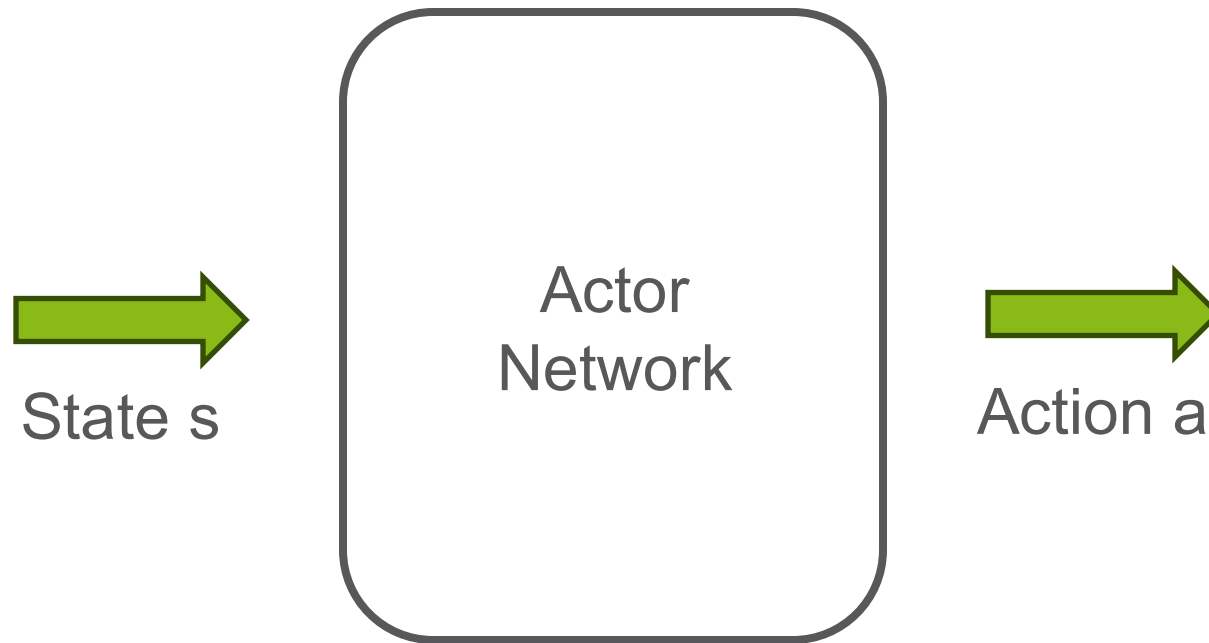$$L(\theta^Q) = \mathbb{E}_{s_t \sim \rho^\beta, a_t \sim \beta, r_t \sim E} \left[ \left( Q(s_t, a_t | \theta^Q) - y_t \right)^2 \right]$$

$$y_t = r(s_t, a_t) + \gamma Q(s_{t+1}, \mu(s_{t+1}) | \theta^Q).$$

State s

Action a

Critic
Network

$Q(s, a)$

# ACTOR

- The actor receives the state as input and produces the appropriate action for the agent.

$$\nabla_{\theta^\mu} J \approx \mathbb{E}_{s_t \sim \rho^\beta} \left[ \nabla_{\theta^\mu} Q(s, a | \theta^Q) |_{s=s_t, a=\mu(s_t | \theta^\mu)} \right]$$

$$= \mathbb{E}_{s_t \sim \rho^\beta} \left[ \nabla_a Q(s, a | \theta^Q) |_{s=s_t, a=\mu(s_t)} \nabla_{\theta_\mu} \mu(s | \theta^\mu) |_{s=s_t} \right]$$
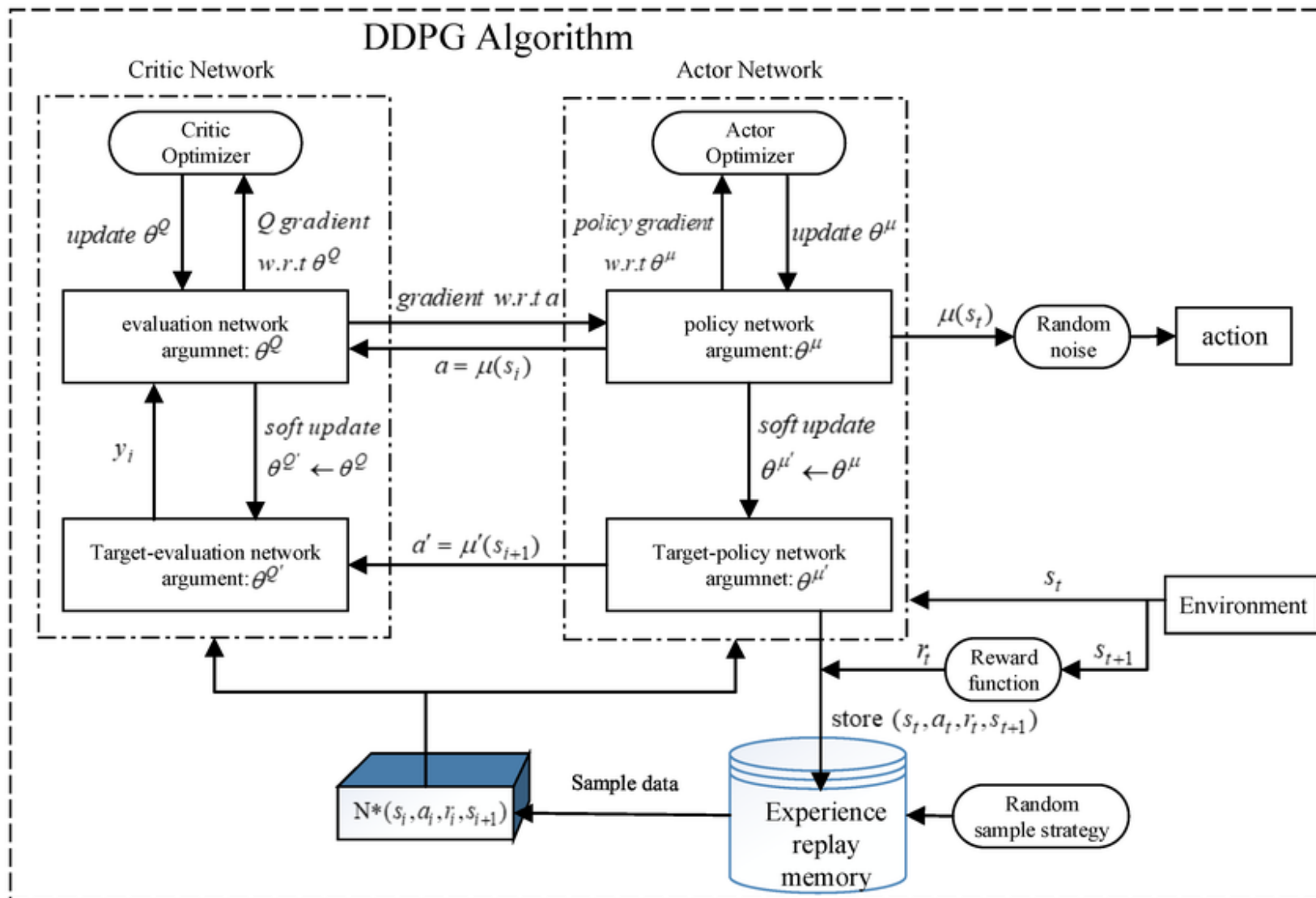


State s — Actor Network — Action a

# IMPLEMENTATION DETAILS

- Due to the deterministic nature of the policy, the agent may fail to explore enough, which can hinder effective training.

- Therefore, a noise component is incorporated into the agent's policy during training, allowing the agent to explore the environment more effectively.

$$a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$$

# DDPG Structure

# DDPG Algorithm

---

**Algorithm 1** DDPG algorithm

---

Randomly initialize critic network $Q(s, a|\theta^Q)$ and actor $\mu(s|\theta^\mu)$ with weights $\theta^Q$ and $\theta^\mu$.

Initialize target network $Q'$ and $\mu'$ with weights $\theta^{Q'} \leftarrow \theta^Q$, $\theta^{\mu'} \leftarrow \theta^\mu$

Initialize replay buffer $R$

**for** episode = 1, M **do**

    Initialize a random process $\mathcal{N}$ for action exploration

    Receive initial observation state $s_1$

    **for** t = 1, T **do**

        Select action $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$ according to the current policy and exploration noise

        Execute action $a_t$ and observe reward $r_t$ and observe new state $s_{t+1}$

        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $R$

        Sample a random minibatch of $N$ transitions $(s_i, a_i, r_i, s_{i+1})$ from $R$

        Set $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$

        Update critic by minimizing the loss: $L = \frac{1}{N}\sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$

        Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N}\sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$

        Update the target networks:

$$\theta^{Q'} \leftarrow \tau\theta^Q + (1-\tau)\theta^{Q'}$$
$$\theta^{\mu'} \leftarrow \tau\theta^\mu + (1-\tau)\theta^{\mu'}$$

    **end for**

**end for**

---