

# AI無線通訊系統實驗

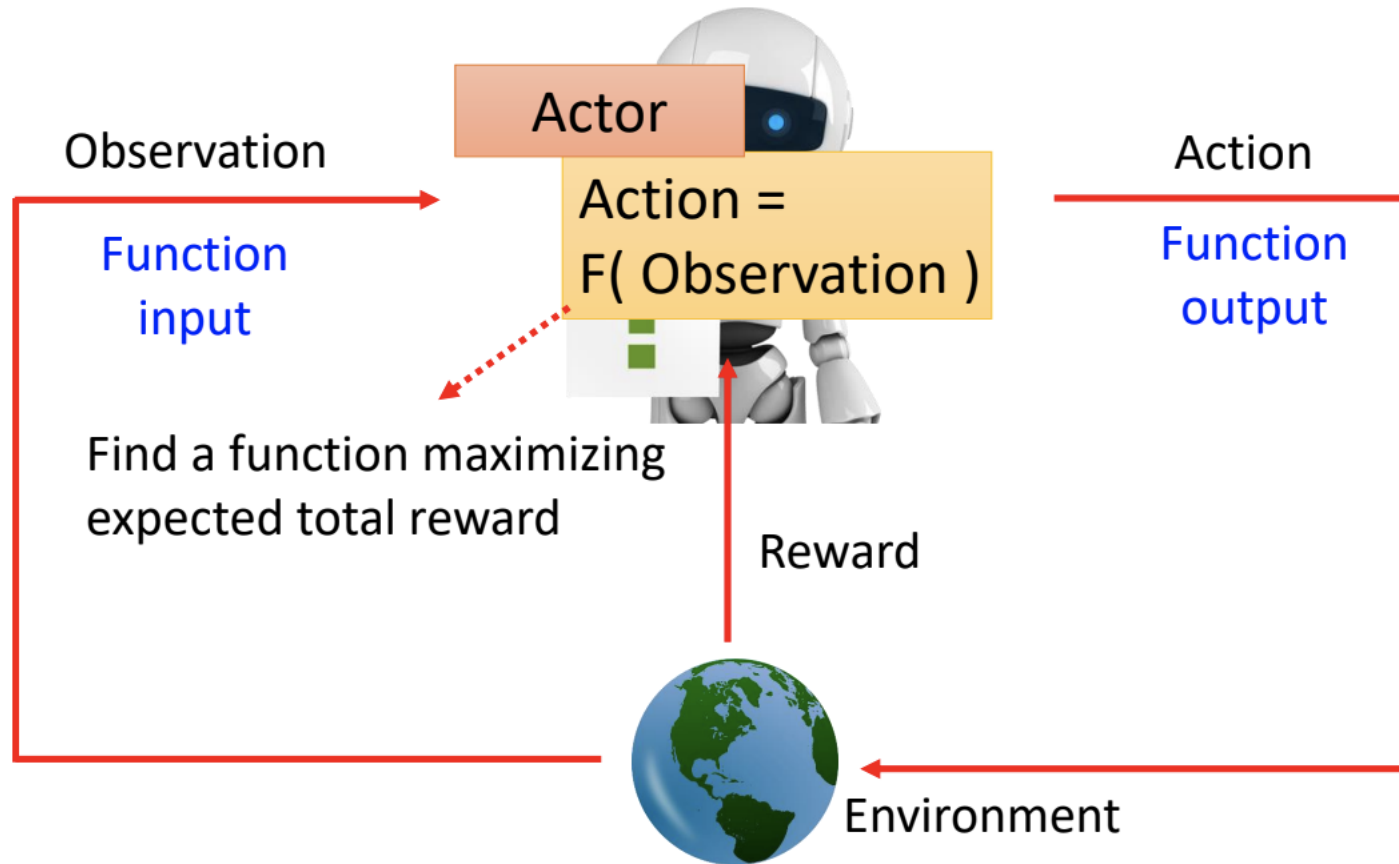
## Reinforcement Learning for Network Resource Allocation

助教: 李育亭、吳淳禎

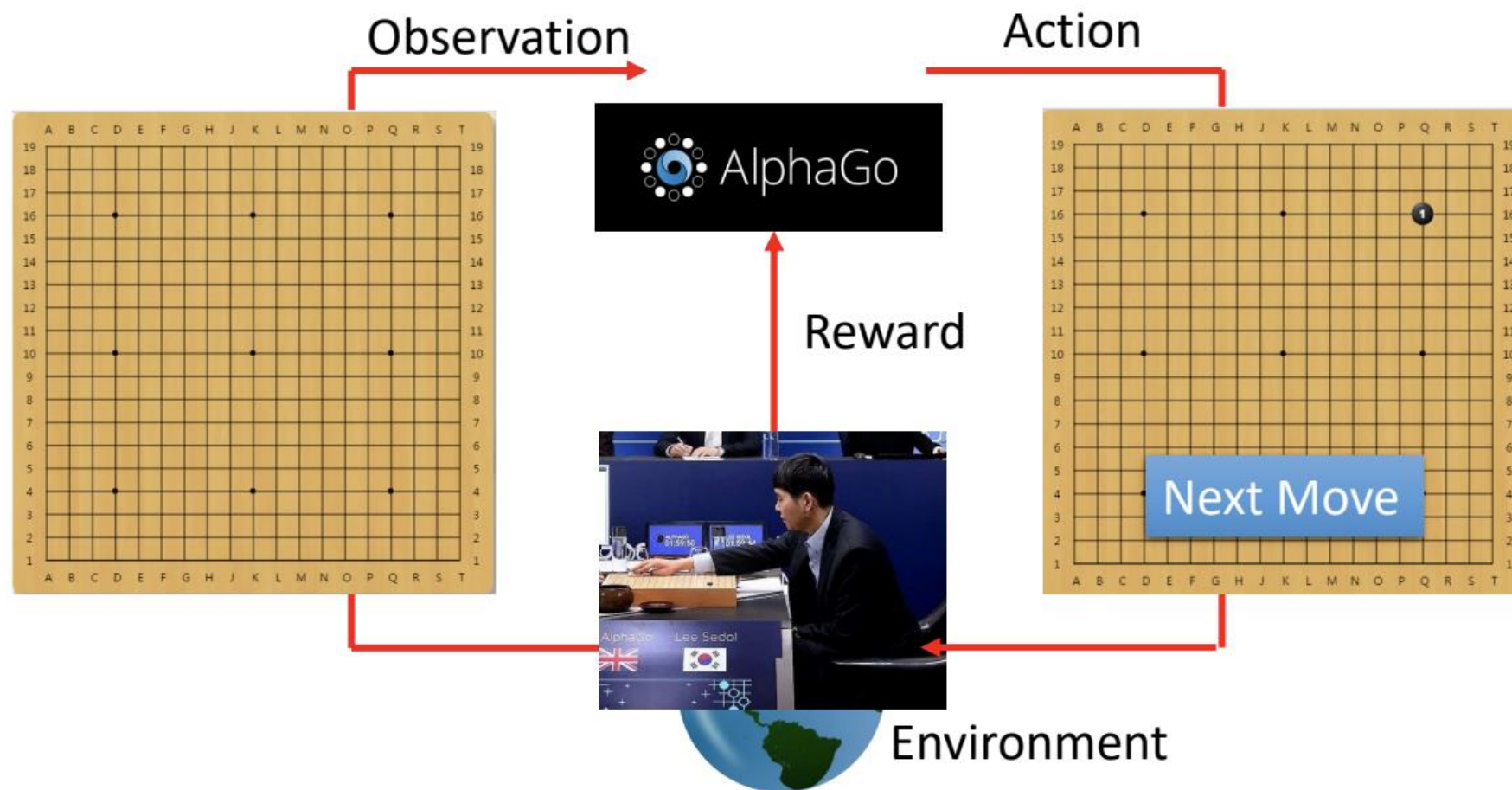
教材編寫: 蕭安紘、沈立翔 博士

# Reinforcement Learning, RL

---



# AlphaGo

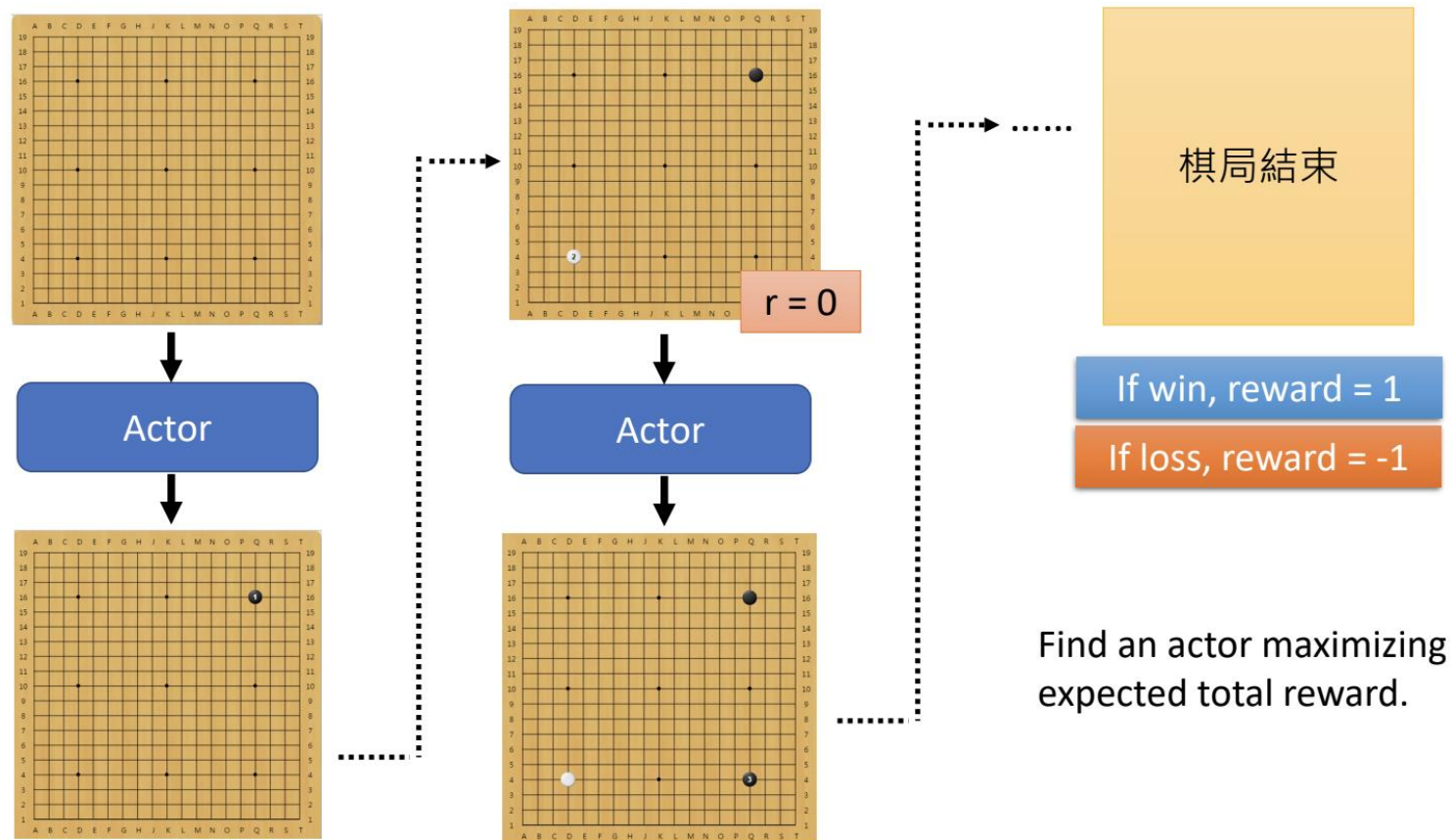


# AlphaGo

Find an actor maximizing expected total reward.



# AlphaGo



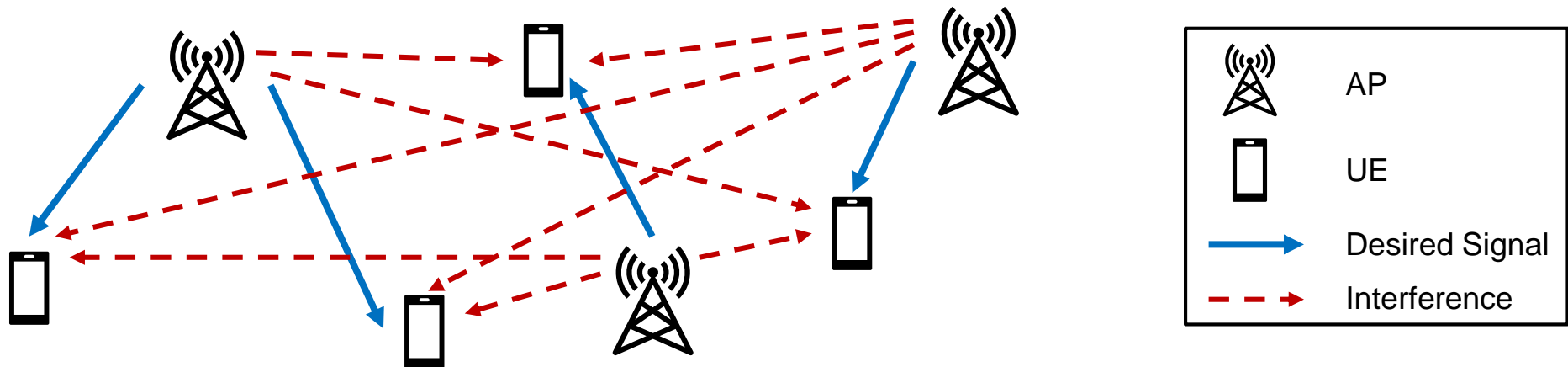
# Resource Allocation

---

- In a real mobile communication, the system always simultaneously performs several tasks under different demanding services from devices
  - However, each device has different channel conditions that we have to take into account for providing optimum resource allocation
- There exist a great number of resources provided from the central system
  - Spatial resources
  - Temporal resource
  - Frequency bandwidth
  - Transmit power

# System Model

- Consider a downlink (DL) scenario with  $N$  access points (APs) and  $K$  user equipments (UEs).



- Channel Gain as the channel gain set from  $N$  APs to  $K$  UEs
$$\mathbf{H} = \{h_{n,k} | \forall 1 \leq n \leq N, 1 \leq k \leq K\}$$
  - $h_{n,k}$ : the channel gain from the  $n$ -th AP to the  $k$ -th UE

# System Model

---

- The downlink transmit power of AP is denoted as

$$\mathbf{P} = \{P_n | \forall n\}$$

- The association indicator set for AP-UE connection

$$\boldsymbol{\rho} = \{\rho_{n,k} \in \{0,1\} | \forall n, k\}$$

- We can obtain the signal-to-noise-ratio (SINR) for the  $n$ -th AP to the  $k$ -th UE as

$$\gamma_{n,k} = \frac{P_n G_{n,k} h_{n,k}}{\sum_{j \neq n}^N P_j G_{j,k} h_{j,k} + \sigma^2}$$

- $G_{n,k}$ : the beamforming gain obtained in the beam training process
  - $\sigma^2$ : background noise
- Then, throughput of the overall system can be expressed as

$$R = \sum_{n=1}^N \sum_{k=1}^K \rho_{n,k} \log_2(1 + \gamma_{n,k})$$



# Problem Formulation

---

- We can then formulate the throughput maximization problem that determines the allocation of power and user association, which is given by

$$\max_{P, \rho} R$$

$$\text{s.t. } 0 < P_n \leq P_{max}, \quad \forall n$$

maximum allowable power constraint

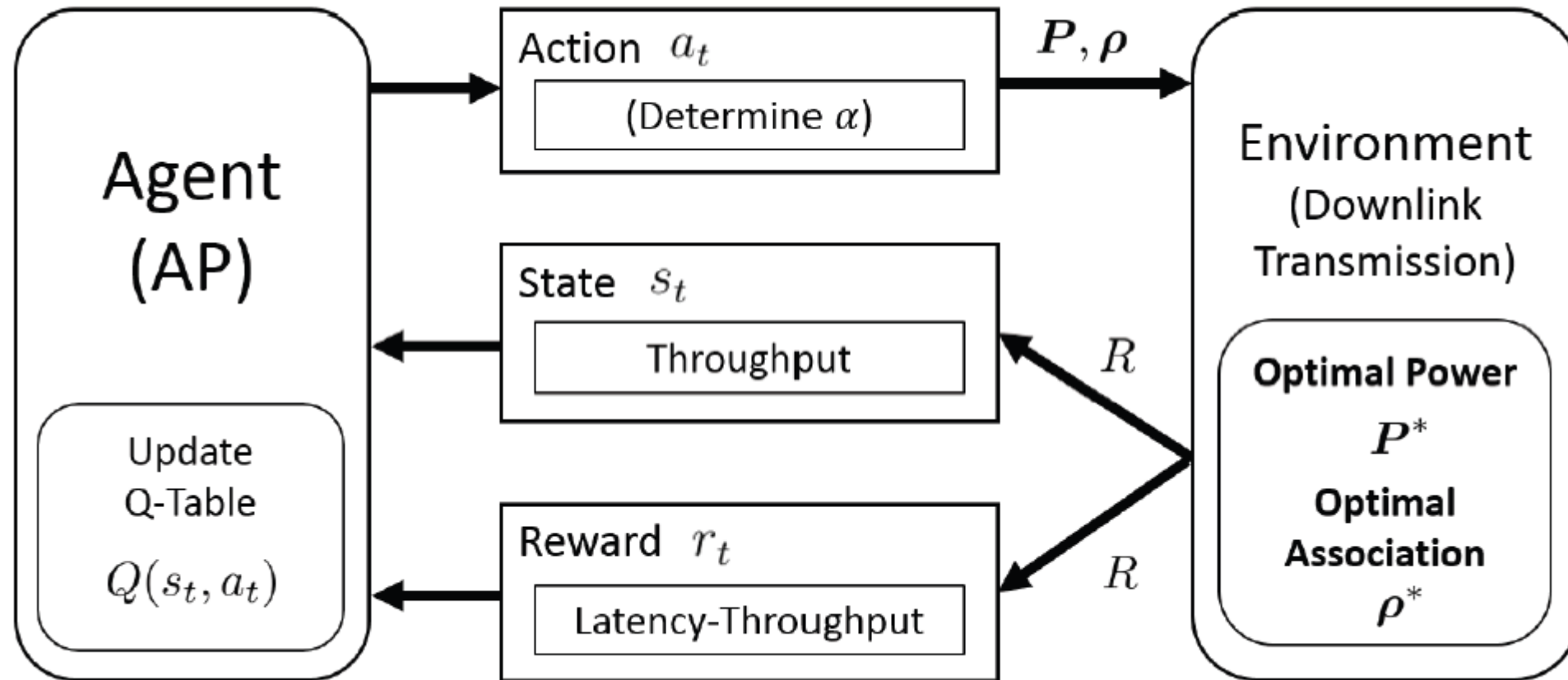
$$\rho_{n,k} \in \{0, 1\}, \quad \forall n, k,$$

association indicator constraint

$$\sum_{n=1}^N \rho_{n,k} \leq 1, \quad \forall k.$$

one UE can only be served by one AP at most

# Q-learning based Resource Allocation



# Q-learning based Resource Allocation

---

- The action is referred as power allocation and association assignment, which is expressed as selection of the maximum value in the Q-table or random action as

$$a_t = \begin{cases} \operatorname{argmax}_{a'_t} Q(s_t, a'_t), & \text{if } \operatorname{rand}() > \epsilon, \\ \text{random action}, & \text{otherwise.} \end{cases}$$

- We exploit  $\epsilon$ -greedy policy, where  $\epsilon \in [0,1]$ .
  - $t$ : the current training or testing time step
- Since the state is a discrete variable, we quantize the sum rate objective as discrete states as

$$s_t = \left\lfloor \frac{R \cdot s_n}{R_{max}} \right\rfloor$$

- $R_{max}$ : the maximum throughput among off-line collected training dataset
  - $s_n$ : pre-defined constant indicating the total number of given states

# Q-learning based Resource Allocation

---

- Rewards

$$r_t = R = \sum_{n=1}^N \sum_{k=1}^K \rho_{n,k} \log_2(1 + \gamma_{n,k})$$

- According to the Bellman equation, the Q-table is updated by

$$Q_t(s_t, a_t) \leftarrow Q_t(s_t, a_t) + \eta \cdot \left[ r_t + \delta \cdot \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) - Q_t(s_t, a_t) \right]$$

- $\eta \in [0,1]$ : learning rate
  - $\delta \in [0,1]$ : discount factor
- In order to address constraint, we often define penalty term to punish the violations.

$$0 < P_n \leq P_{max}, \quad \forall n$$

# Q-learning based Resource Allocation

---

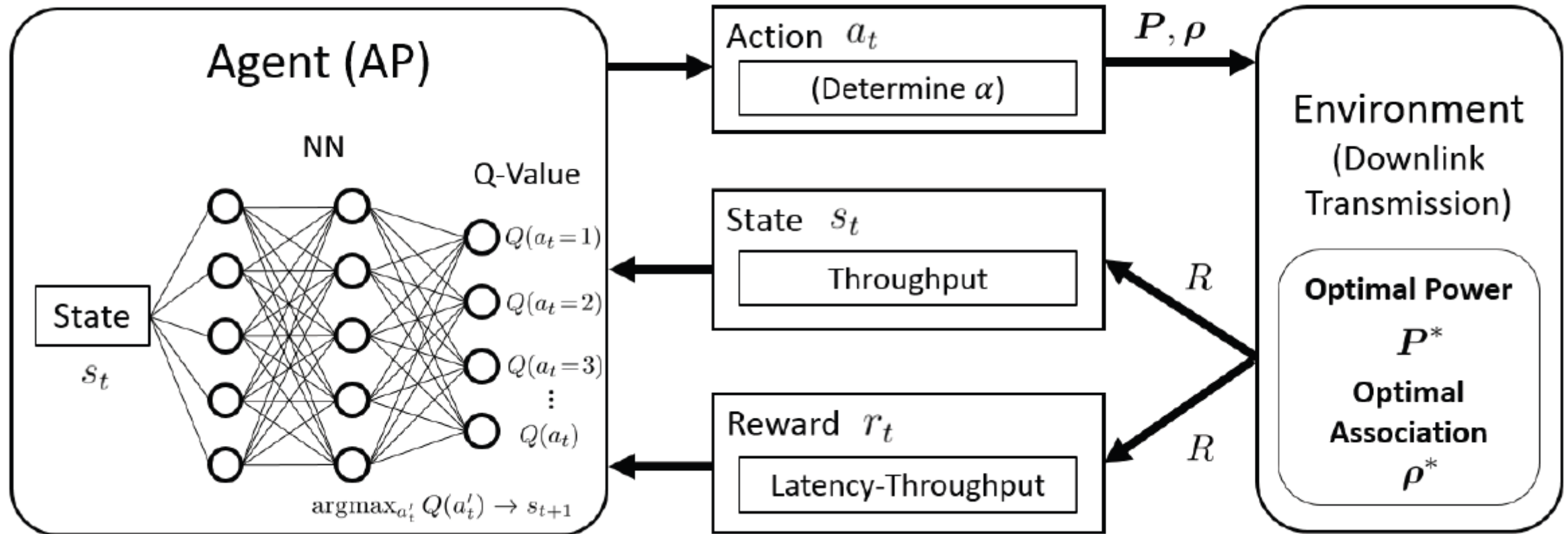
---

**Algorithm 8: Q-Learning for Resource Allocation**

---

- 1: **Initialization:** Time step  $t = 0$ , learning rate  $\eta$ , discount factor  $\delta$ , exploitation probability  $\epsilon$
  - 2: **Offline:** Collect training dataset of power and association, find  $R_{max}$ , and update the Q-table based on **Q-learning**
  - 3: **Online:** Based on trained Q-table, perform online testing based on **Q-learning**
  - 4: **Q-learning:**
  - 5: **while** Not converged as  $t < T$  **do**
  - 6:    $t = t + 1$
  - 7:   Randomly generate a number  $\kappa \in [0, 1]$
  - 8:   **if**  $t = 1$  or  $\kappa > \epsilon$  **then**
  - 9:     Select random action  $a_t$
  - 10:  **else**
  - 11:   Select optimal action based on  $a_t = \underset{a'_t}{\operatorname{argmax}} Q(s_t, a'_t)$
  - 12:  **end if**
  - 13:   Obtain power  $\mathbf{P}$  and association  $\rho$  decision based on  $a_t$
  - 14:   Obtain reward:  $r_t = R$
  - 15:   Update Q-table:  $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \eta \cdot [r_t + \delta \cdot \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$
  - 16:   Update next state:  $s_{t+1} = \lceil \frac{r_t \cdot s_n}{R_{max}} \rceil$
  - 17: **end while**
-

# Deep Q-learning



# Deep Q-learning

