**Value-based RL：**
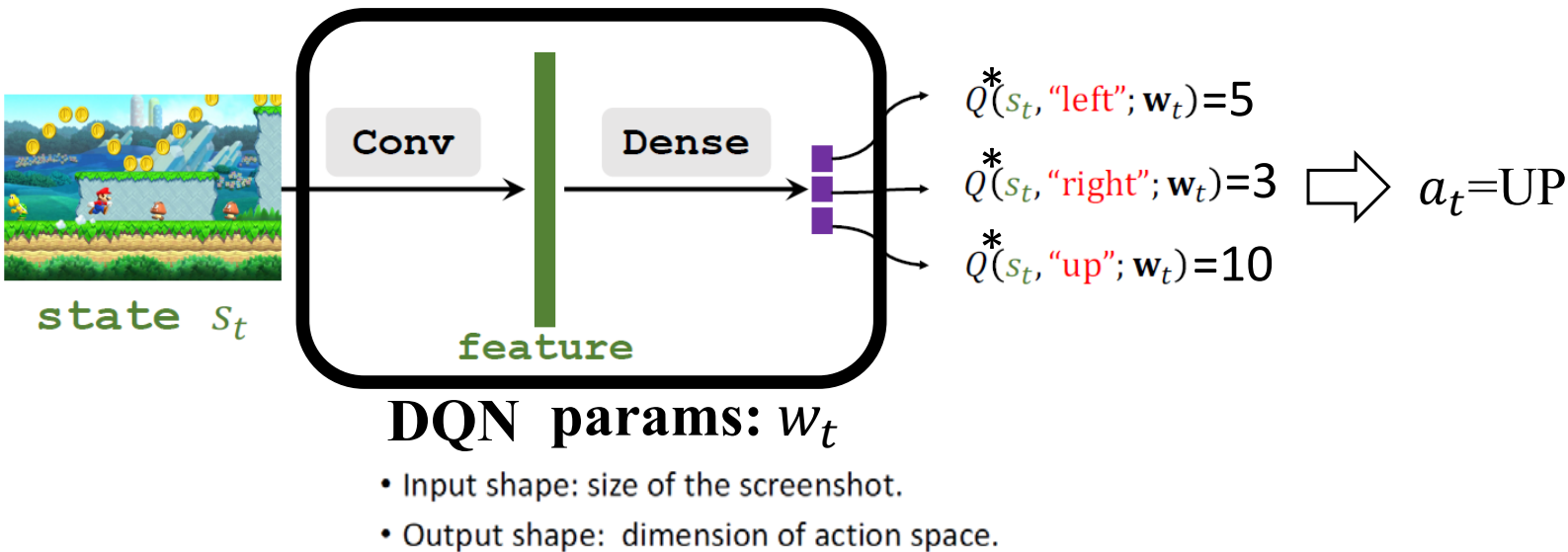**用神经网络去拟合optimal action-value function Q***

**Goal:** Win the game ($\approx$ maximize the total reward.)

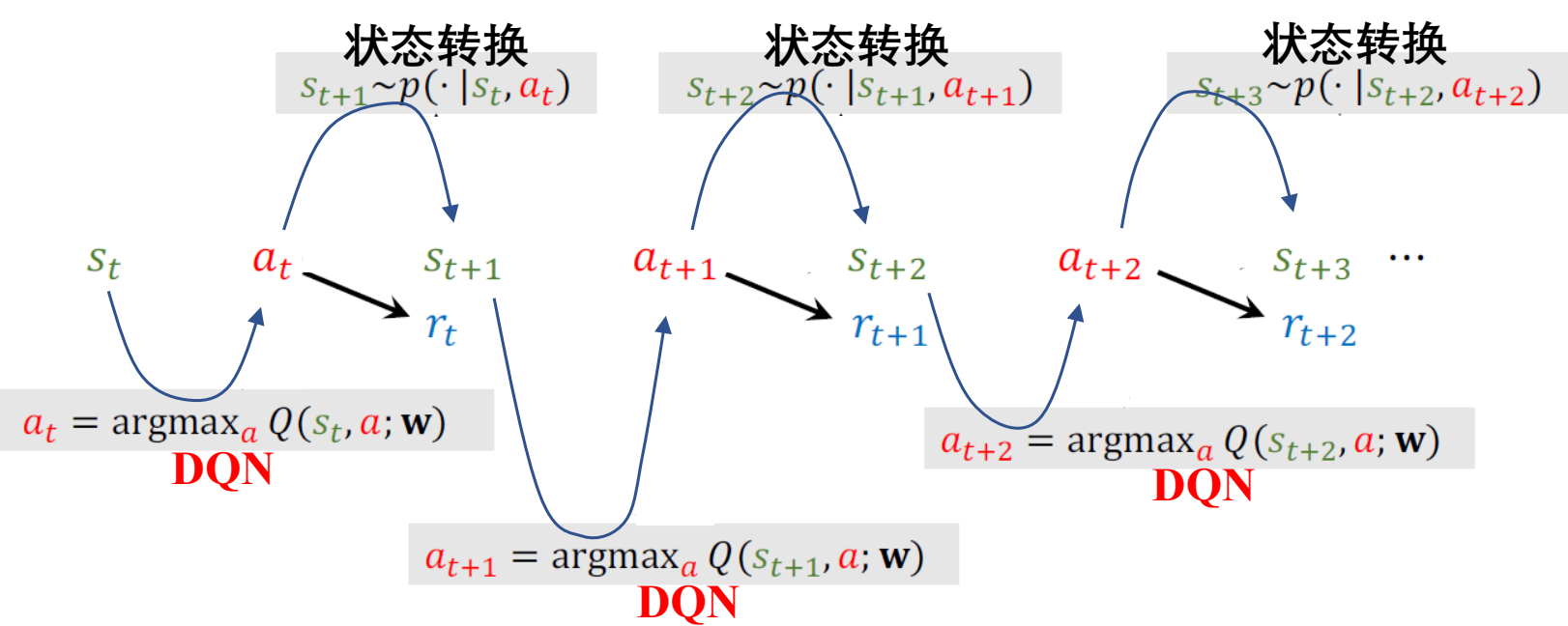**Question:** If we know $Q^\star(s, a)$, what is the best action?

- Obviously, the best action is $a^\star = \operatorname*{argmax}\limits_{a} Q^\star(s, a)$.

**Challenge:** We do not know $Q^\star(s, a)$.

- Solution: Deep Q Network (DQN)
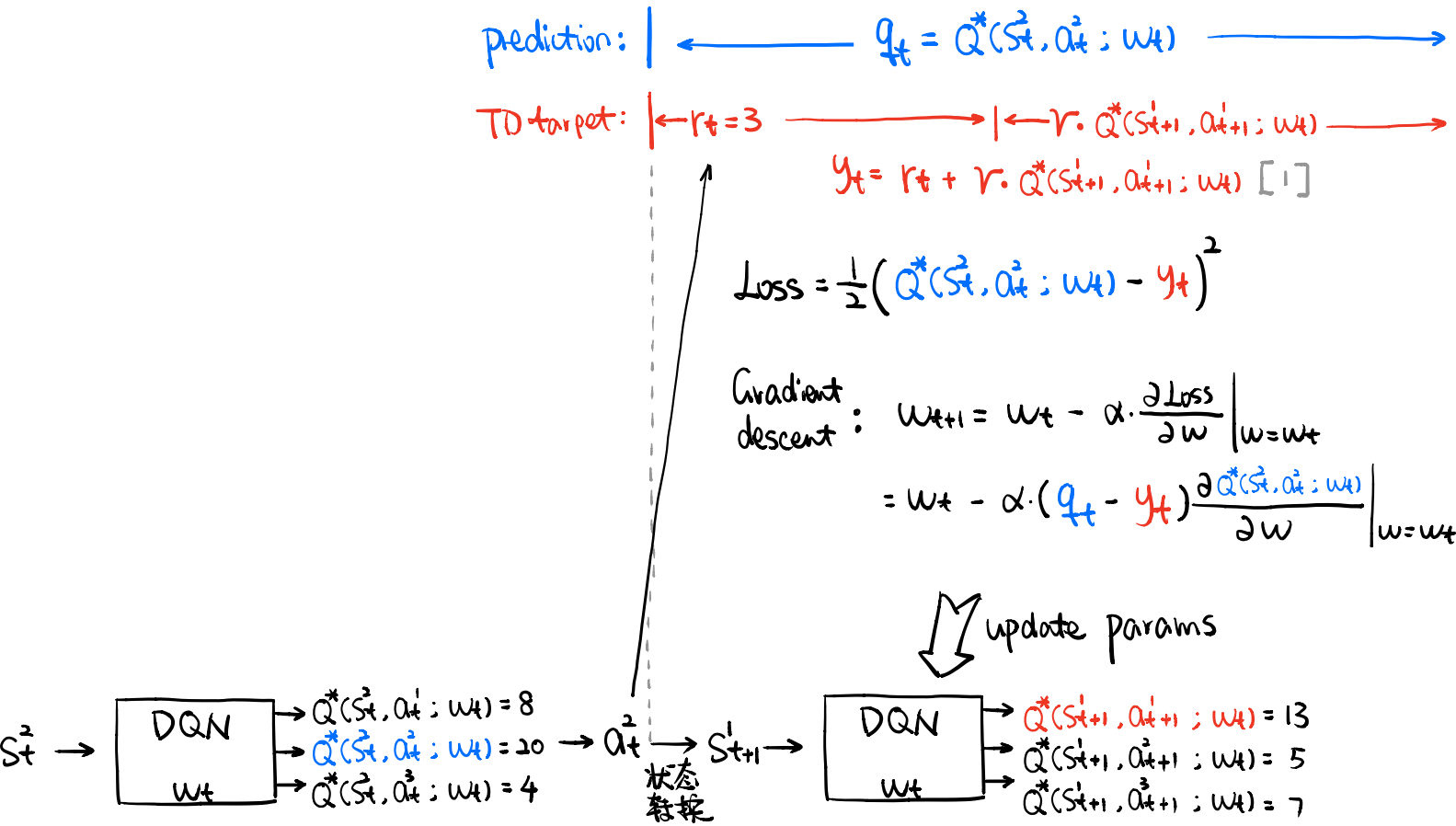- Use neural network $Q(s, a; \mathbf{w})$ to approximate $Q^\star(s, a)$



$Q^*(s_t, \text{"left"}; \mathbf{w}_t) = 5$
$Q^*(s_t, \text{"right"}; \mathbf{w}_t) = 3 \Rightarrow a_t = \text{UP}$
$Q^*(s_t, \text{"up"}; \mathbf{w}_t) = 10$

**state** $s_t$

**DQN params:** $w_t$

- Input shape: size of the screenshot.
- Output shape: dimension of action space.

**trained DQN:** 假如我们有一个trained DQN (Q*)，如何将其应用于agent?



状态转换
$s_{t+1} \sim p(\cdot \mid s_t, a_t)$

状态转换
$s_{t+2} \sim p(\cdot \mid s_{t+1}, a_{t+1})$

状态转换
$s_{t+3} \sim p(\cdot \mid s_{t+2}, a_{t+2})$

$s_t \quad a_t \quad s_{t+1} \quad a_{t+1} \quad s_{t+2} \quad a_{t+2} \quad s_{t+3} \cdots$
$r_t \quad\quad r_{t+1} \quad\quad r_{t+2}$

$a_t = \operatorname*{argmax}\limits_{a} Q(s_t, a; \mathbf{w})$
**DQN**

$a_{t+2} = \operatorname*{argmax}\limits_{a} Q(s_{t+2}, a; \mathbf{w})$
**DQN**

$a_{t+1} = \operatorname*{argmax}\limits_{a} Q(s_{t+1}, a; \mathbf{w})$
**DQN**

**问:** 我们要利用什么数据，如何训练一个DQN来去拟合Q*?

# Temporal Difference (TD) Learning

Prediction: $\longleftarrow \quad q_t = Q^*(s_t^2, a_t^2; w_t) \longrightarrow$

TD target: $\leftarrow r_t = 3 \longrightarrow \mid \leftarrow \gamma \cdot Q^*(s_{t+1}^1, a_{t+1}^1; w_t) \longrightarrow$

$$y_t = r_t + \gamma \cdot Q^*(s_{t+1}^1, a_{t+1}^1; w_t) \quad [1]$$

$$\text{Loss} = \frac{1}{2}\left(Q^*(s_t^2, a_t^2; w_t) - y_t\right)^2$$

Gradient descent: $\quad w_{t+1} = w_t - \alpha \cdot \dfrac{\partial \text{Loss}}{\partial w}\bigg|_{w=w_t}$

$$= w_t - \alpha \cdot (q_t - y_t)\dfrac{\partial Q^*(s_t^2, a_t^2; w_t)}{\partial w}\bigg|_{w=w_t}$$

⇩ update params

$S_t^2 \rightarrow$ | DQN  $w_t$ | $\rightarrow$ $Q^*(s_t^2, a_t^1; w_t) = 8$
$Q^*(s_t^2, a_t^2; w_t) = 20$
$Q^*(s_t^2, a_t^3; w_t) = 4$
$\rightarrow a_t^2 \rightarrow$ 状态转移 $\rightarrow S_{t+1}^1 \rightarrow$ | DQN  $w_t$ | $\rightarrow$ $Q^*(s_{t+1}^1, a_{t+1}^1; w_t) = 13$
$Q^*(s_{t+1}^1, a_{t+1}^2; w_t) = 5$
$Q^*(s_{t+1}^1, a_{t+1}^3; w_t) = 7$

## THEN:

$S_{t+1}^1 \rightarrow$ | DQN  $w_{t+1}$ | $\rightarrow$ $Q^*(s_{t+1}^1, a_{t+1}^1; w_t)$
$Q^*(s_{t+1}^1, a_{t+1}^2; w_t)$
$Q^*(s_{t+1}^1, a_{t+1}^3; w_t)$
$\longrightarrow$ 重复上述过程，可以看出，在每个时间点上，就会利用 TD learning 更新一次权重 $w$。

[1]:

$$Q^*(s_t^2, a_t^2; w_t) = \max_{\pi} Q_\pi(s_t^2, a_t^2; w_t)$$
$$= \max_{\pi} E[\,U_t \mid S_t = s_t^2, A_t = a_t^2 ; w_t\,]$$

$$Q^*(s_{t+1}^1, a_{t+1}^1; w_t) = \max_{\pi} Q_\pi(s_{t+1}^1, a_{t+1}^1; w_t)$$
$$= \max_{\pi} E[\,U_{t+1} \mid S_{t+1} = s_{t+1}^1, A_{t+1} = a_{t+1}^1 ; w_t\,]$$

> **Identity:** $U_t = R_t + \gamma \cdot U_{t+1}.$

- $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot R_{t+3} + \gamma^4 \cdot R_{t+4} + \cdots$
  $= R_t + \gamma \cdot \underbrace{(R_{t+1} + \gamma \cdot R_{t+2} + \gamma^2 \cdot R_{t+3} + \gamma^3 \cdot R_{t+4} + \cdots)}_{= U_{t+1}}$

$$\underbrace{r_t + \gamma \cdot Q^*(s_{t+1}^1, a_{t+1}^1; w_t) = y_t}_{\text{Target}} \approx \underbrace{Q^*(s_t^2, a_t^2; w_t)}_{\text{Prediction}}$$

# Temporal Difference (TD) Learning

**Algorithm:** One iteration of TD learning.

1. Observe state $S_t = s_t$ and perform action $A_t = a_t$.

2. Predict the value: $q_t = Q(s_t, a_t; \mathbf{w}_t)$.

3. Differentiate the value network: $\mathbf{d}_t = \frac{\partial\, Q(s_t, a_t; \mathbf{w})}{\partial\, \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.

4. Environment provides new state $s_{t+1}$ and reward $r_t$.

5. Compute TD target: $y_t = r_t + \gamma \cdot \max_{a} Q(s_{t+1}, a; \mathbf{w}_t)$.

6. Gradient descent: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot (q_t - y_t) \cdot \mathbf{d}_t$.

$$S_{t+1} \rightarrow \boxed{DQN} \equiv \frac{\underline{\;Q^*(S_{t+1}, a_1; w_t)}}{\vdots}$$