

**Blockchain at University of Batna 2**

**Reinforcement Learning**

**Agents For Finance**

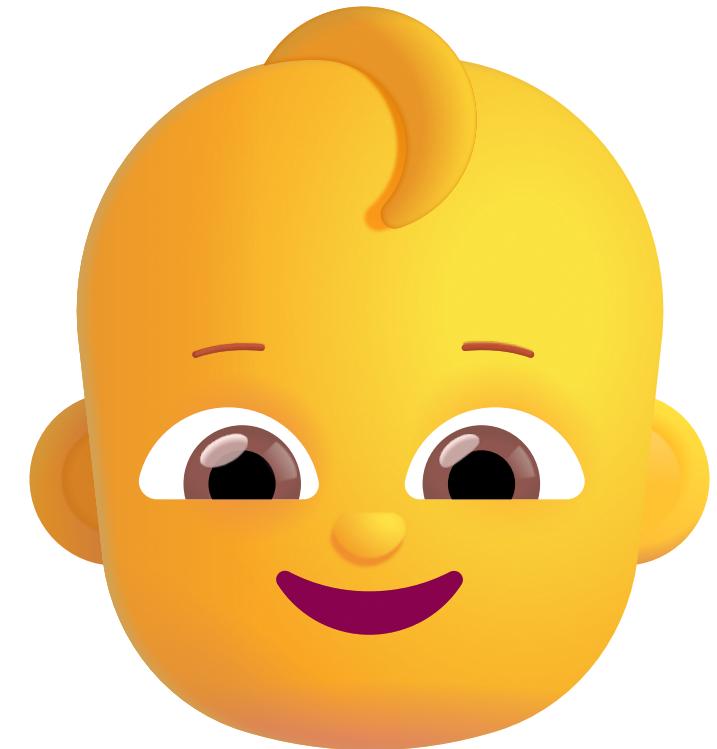
**By Kamel Brouthen**



# Kamel Brouthen

- ESI Algiers CS Student
- GDG Algiers Dev Core Team Member
- SOAI Algiers Technical Manager

# **Building Intuition**

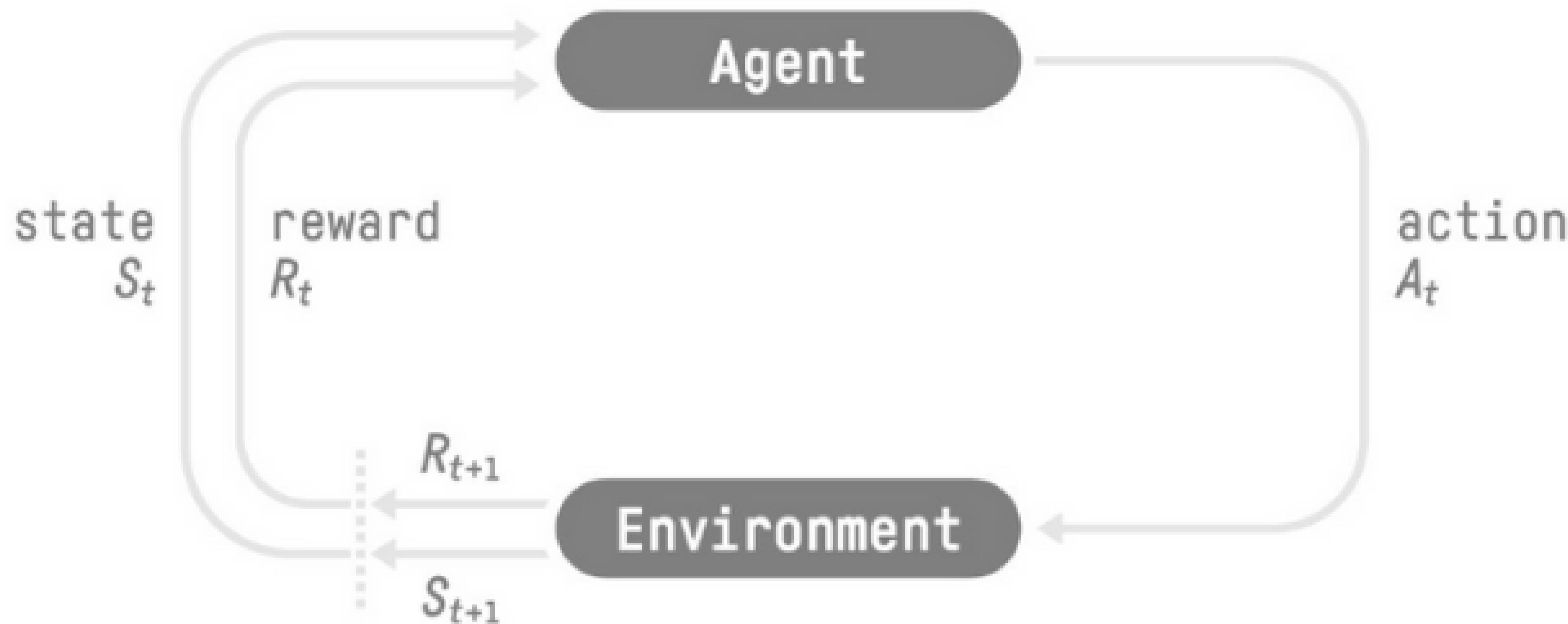


**How did we learn about the world?**

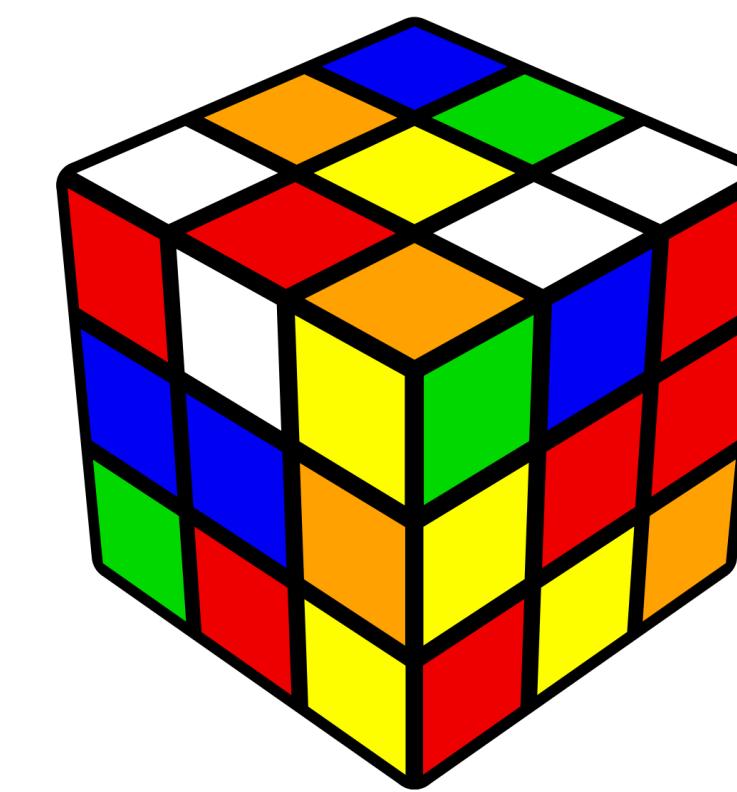
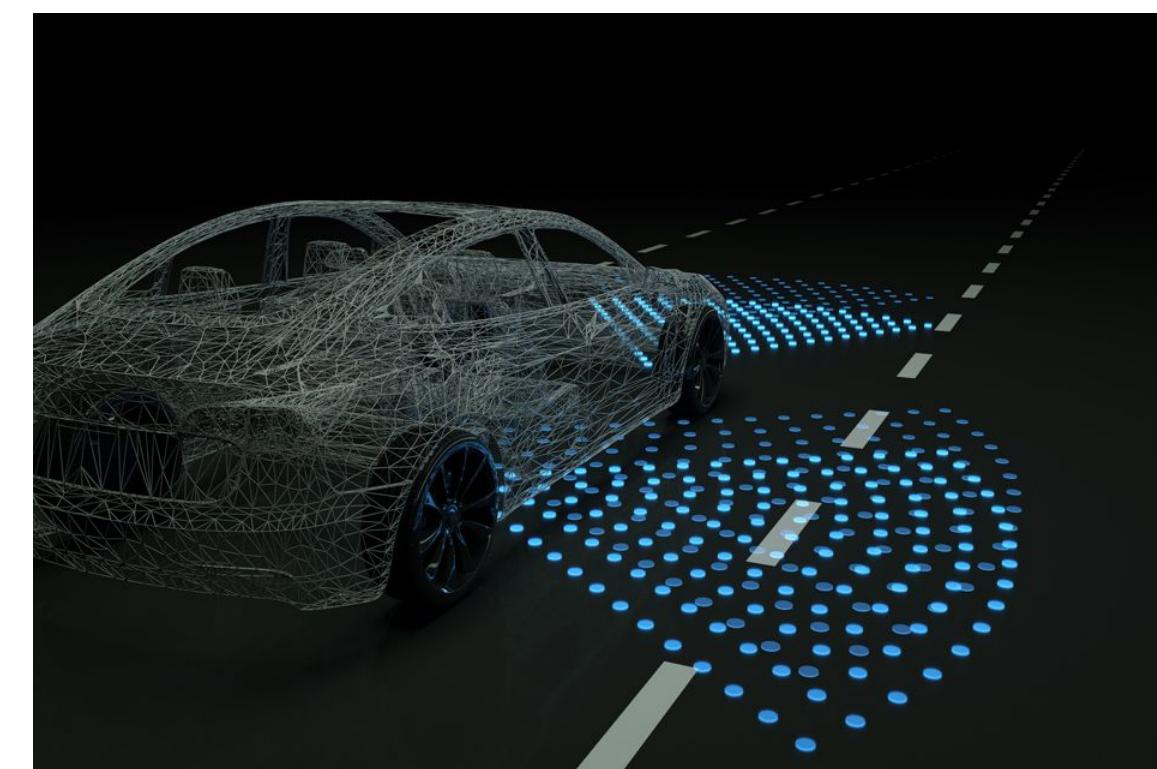
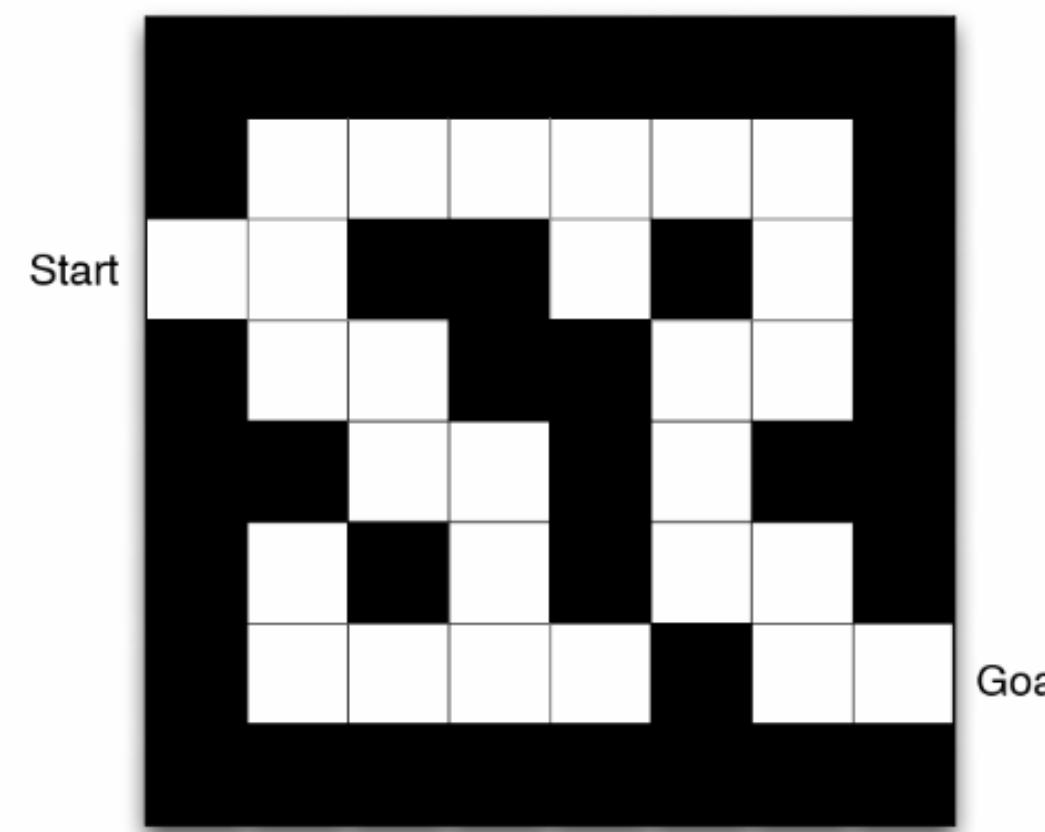
# **Reinforcement Learning**

The AI discipline of learning through **interaction** and **feedback** signal from a given **environment**  
**(Trial & Error)**

# Reinforcement Learning Framework



# Reinforcement Learning Illustrated



# Reinforcement Learning (Maybe)



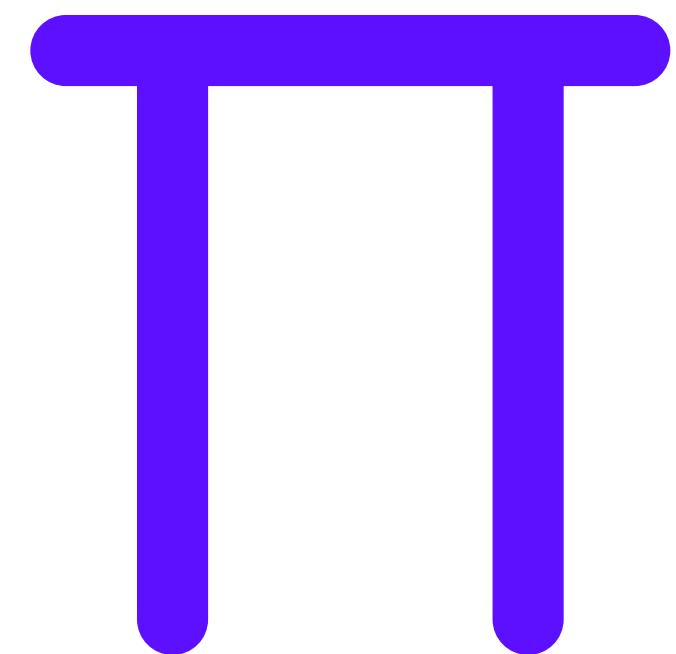
# The Reward Hypothesis (why RL works)

“All of what we mean by **goals** and purposes can be well thought of as **maximization** of the **expected value** of the **cumulative sum** of a received scalar signal (**reward**).”

Sutton, R. S. The reward hypothesis. 2004

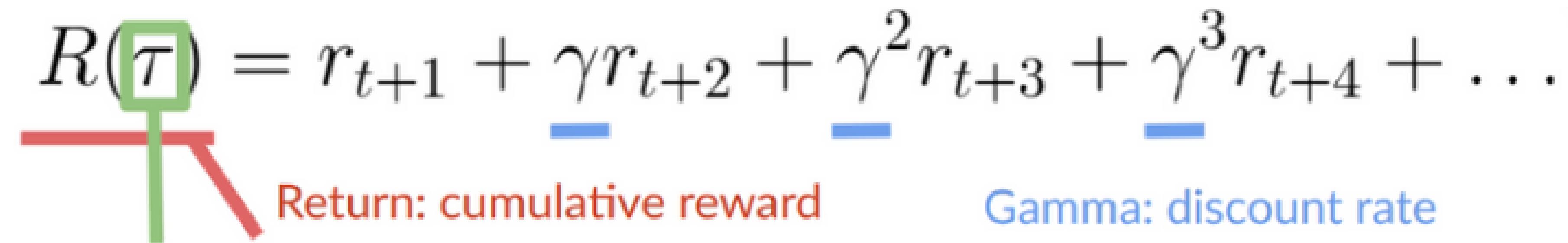
<http://incompleteideas.net/rli.cs.ualberta.ca/RLAI/rewardhypothesis.html>

# Mathematics (AKA AI)



The building blocks of solving RL problems

# Return (Cumulative Reward)

$$R(\tau) = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots$$


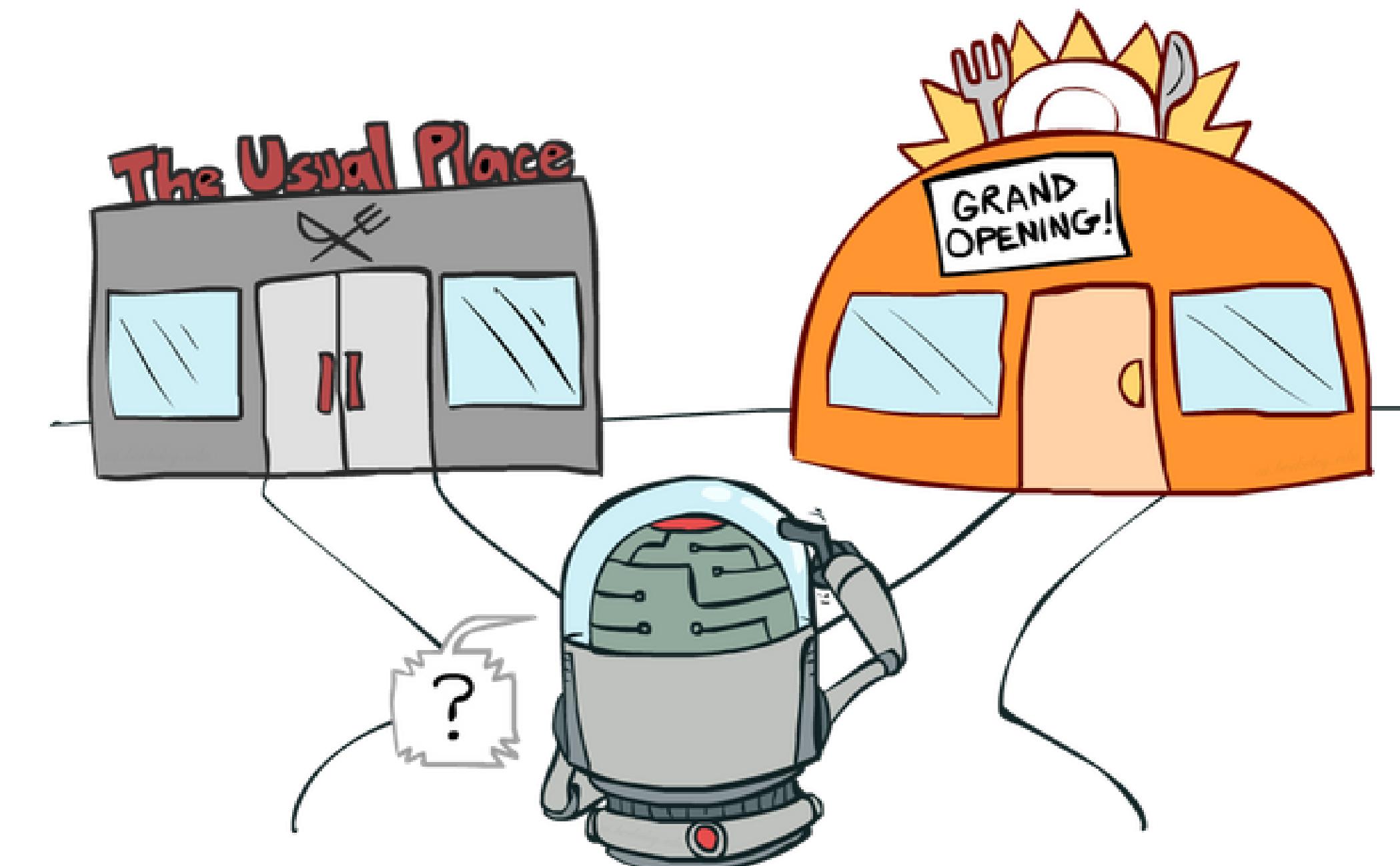
Trajectory (read Tau)  
Sequence of states and actions

Gamma: discount rate

**Goal:** Maximize the expected return

# Exploration VS Exploitation

(to keep in mind)





# Solving RL

## Value Based Methods

**Estimating** how **good** is it  
to be in a given **state** and  
taking a particular **action**

## Policy Based Methods

**Directly** adjusting the  
agent's **policy** over the  
**actions** space

# Value Based Method (Q-Values)

$$Q^\pi(s, a) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right]$$

**Explanation:** Being at **state s** and taking **action a**,  
how much **return** do we **expect**

# Temporal Difference Learning

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

New Q-value estimation      Former Q-value estimation      Learning Rate      Immediate Reward      Discounted Estimate optimal Q-value of next state      Former Q-value estimation

TD Target

TD Error

**Explanation:** Q-Values estimations get adjusted through experience by the immediate reward and bootstrapping on next trajectories

# The Policy (given final estimations)

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

**Explanation:** For each state  $s$ , we take the action  $a$  that maximizes the state-action value

# Building Intuition

## (Mouse wants Cheese)



# Q-Learning Algorithm

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  $Q(\text{terminal}, \cdot) = 0$

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  $R, S'$

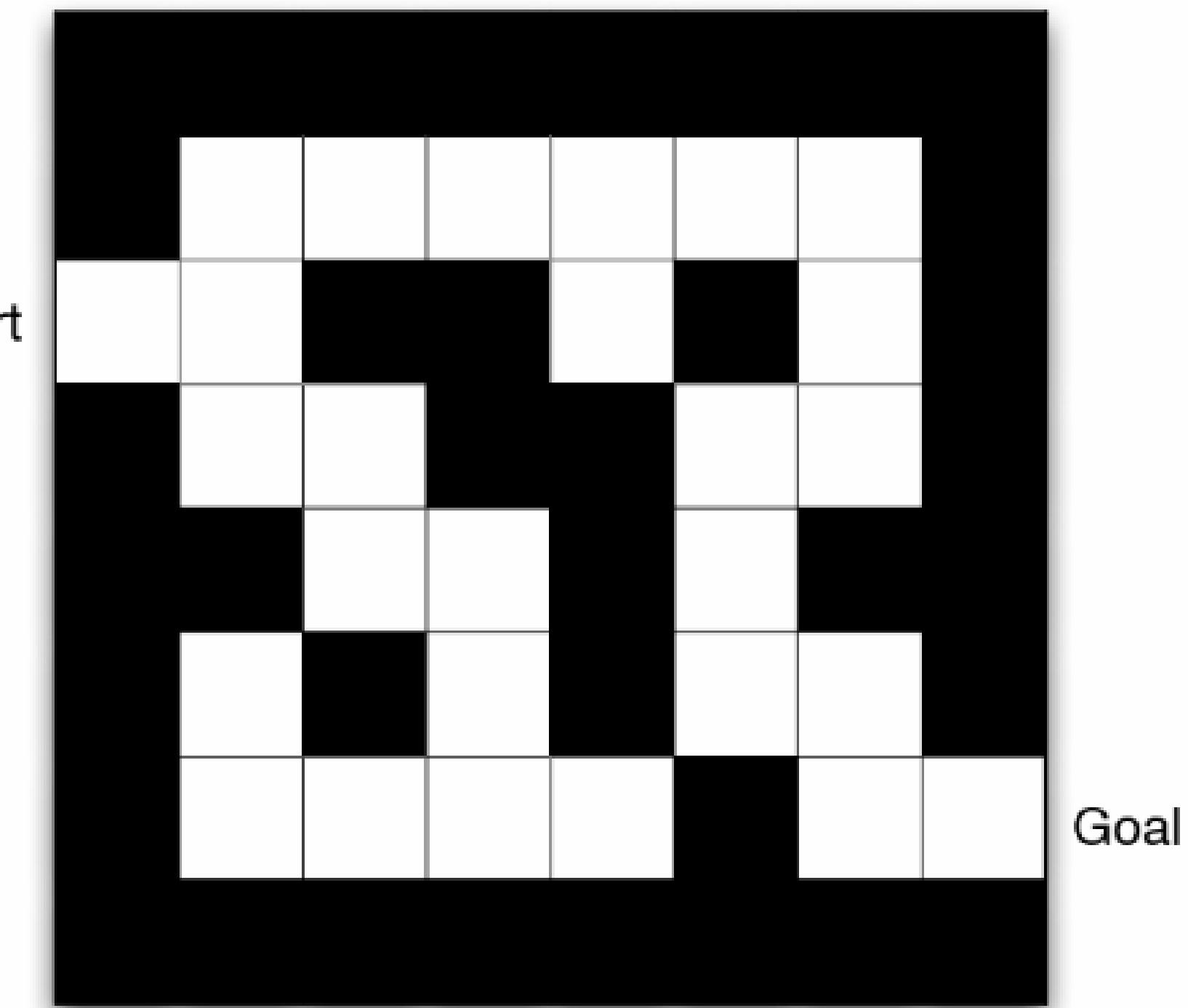
$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

$S \leftarrow S'$

    until  $S$  is terminal

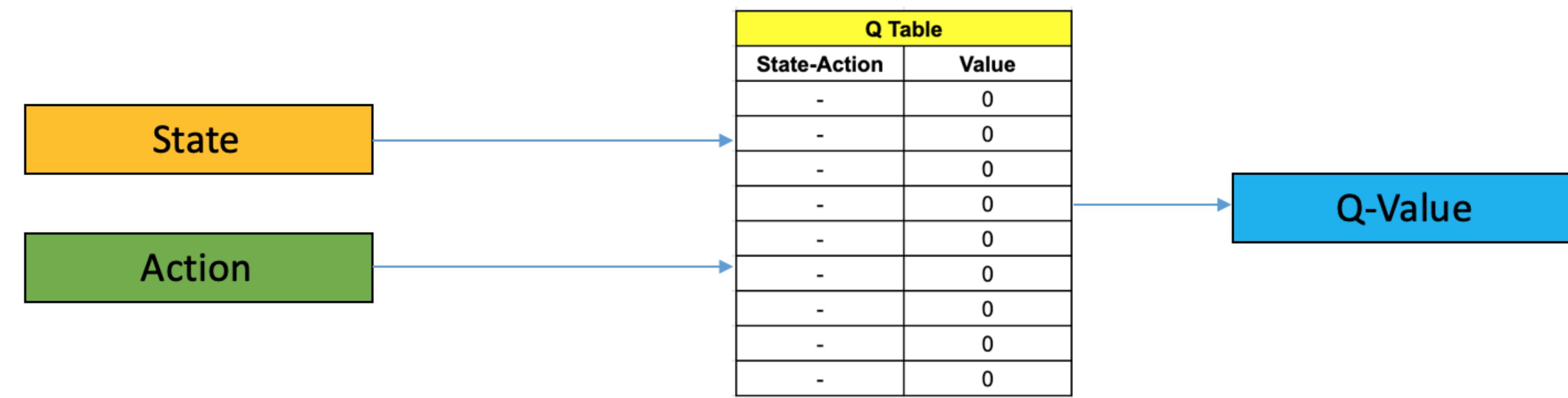
[https://medium.com/analytics-vidhya/q-learning-expected-sarsa-and-comparison-of-td-learning-algorithms-e4612064de97](https://medium.com.analytics-vidhya/q-learning-expected-sarsa-and-comparison-of-td-learning-algorithms-e4612064de97)

# Applied Q-Learning (Navigate the Maze)

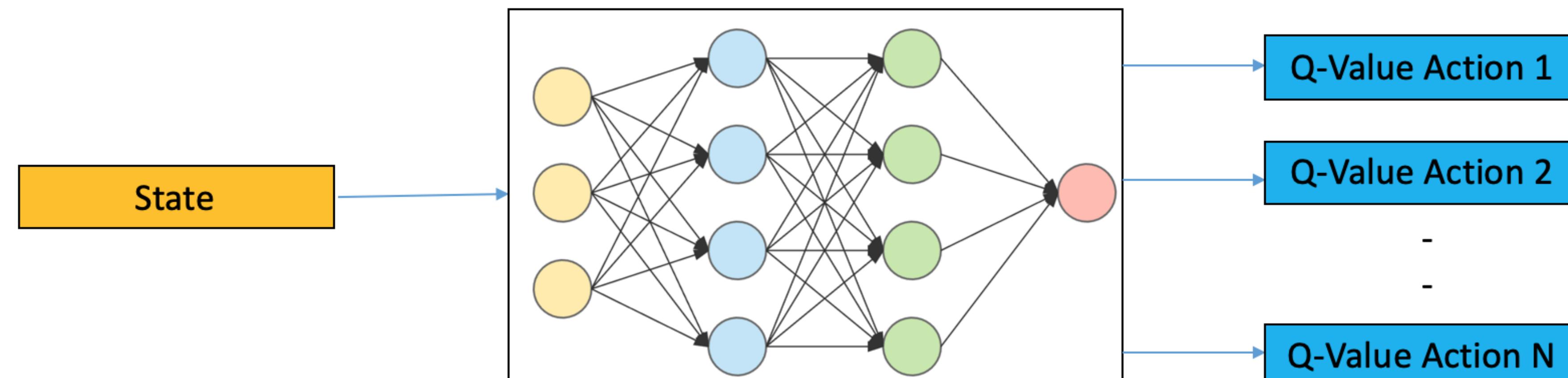


# Q-Learning Limitations

- Handling high **dimensional** state/action spaces
- Function approximation **generalization** capabilities
- Familiarity and usability only in **discrete** action spaces



## Q Learning

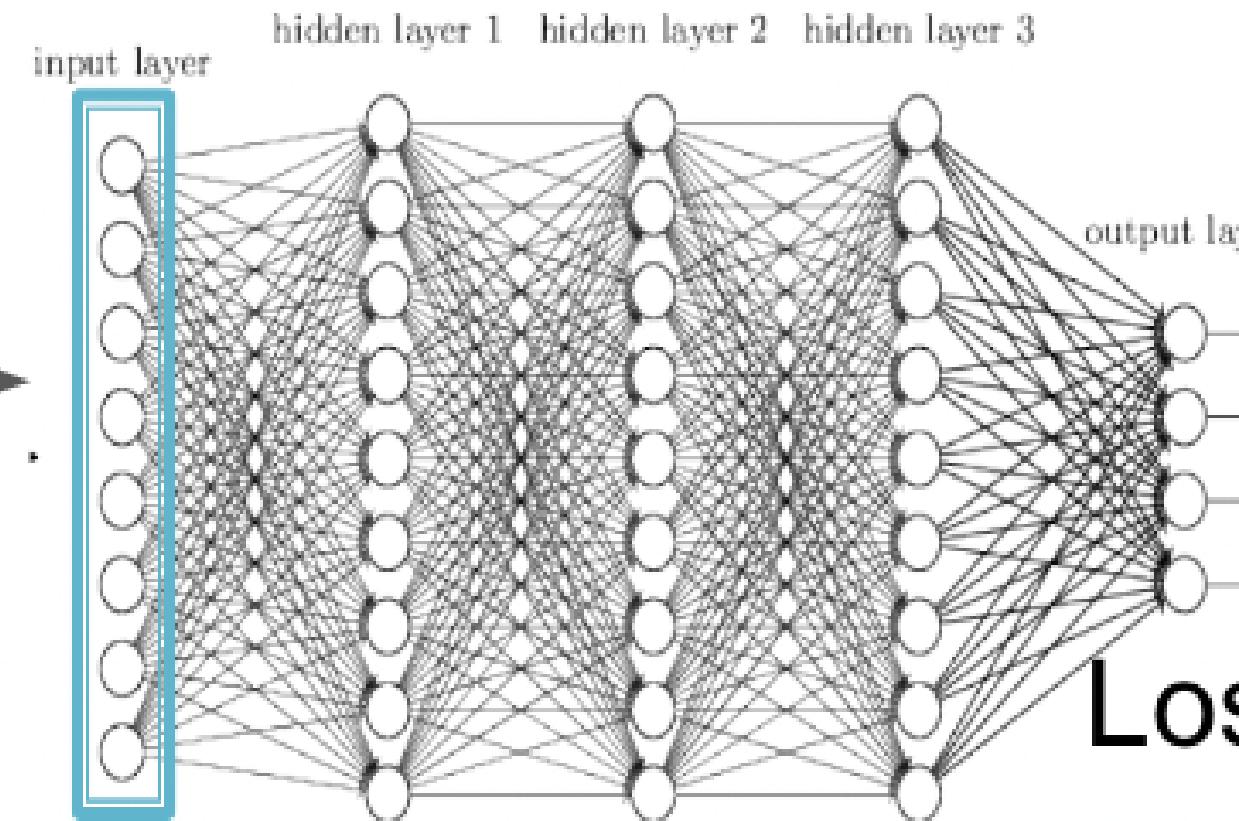


## Deep Q Learning

# Deep Q-Learning

**Q-Network**

State-S →



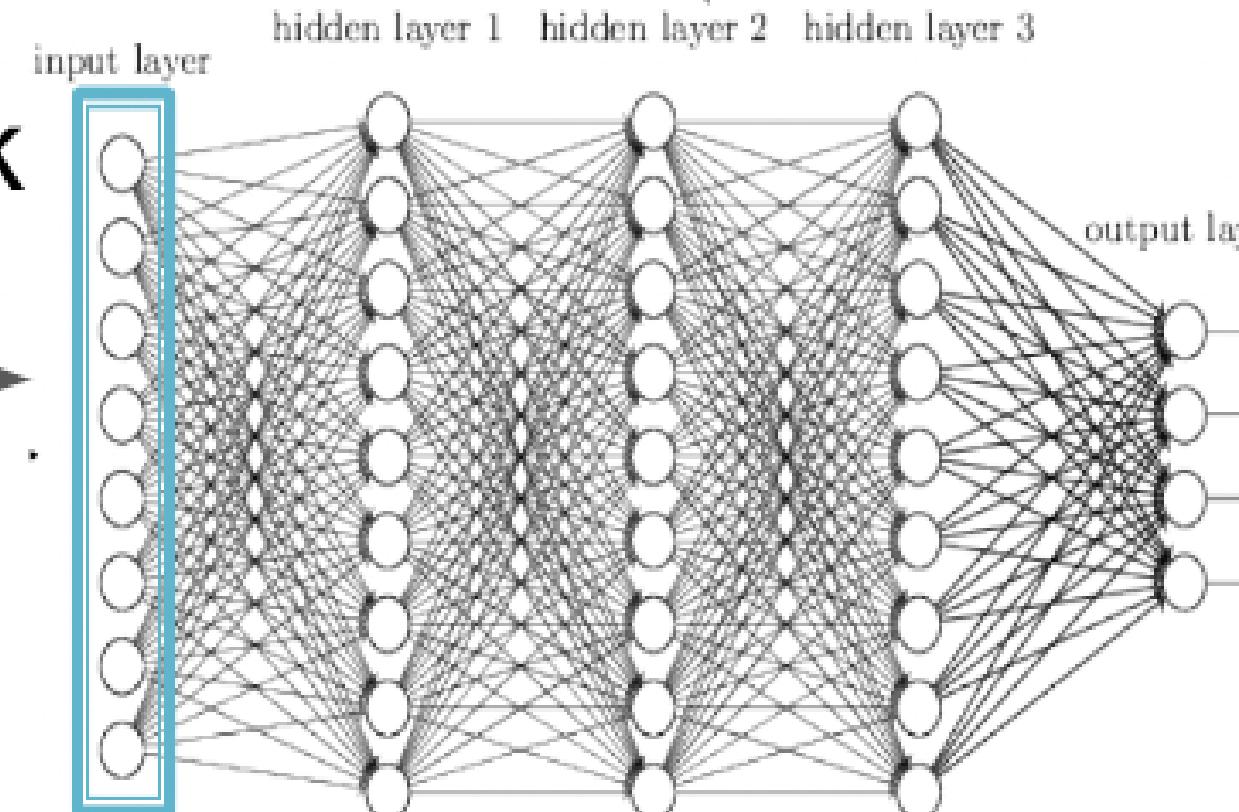
$$Q(s, a; \theta)$$

$$\text{Loss} = \left( r + \gamma \max_{a'} Q(s', a'; \theta^-_i) - Q(s, a; \theta_i) \right)^2$$

↓  
θ updates θ<sup>-</sup>  
every C timesteps

**Target-Network**

State-S →



$$Q_T(s', a'; \theta^-)$$

# Deep Q-Learning

---

## 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  $\mathcal{D}$

        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**

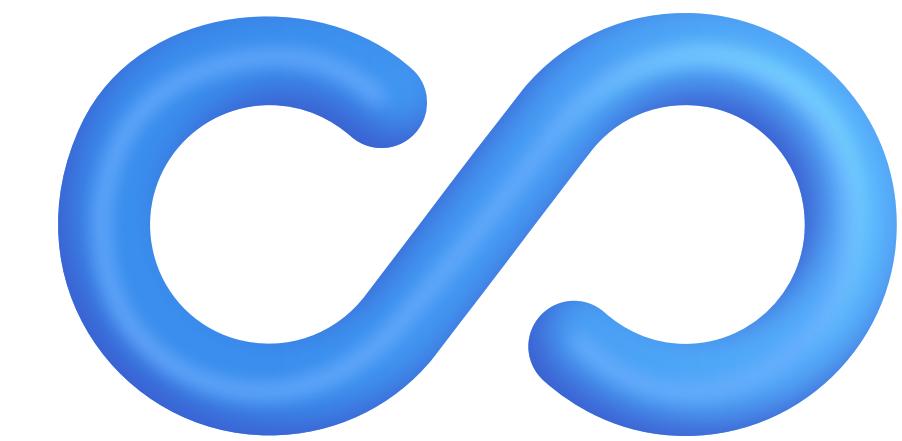
**end for**

---

# Applied Deep Q-Learning (Bitcoin Trading)



# **Human Learners!**



**Agents are learning, so should we...**



**Thank you!**