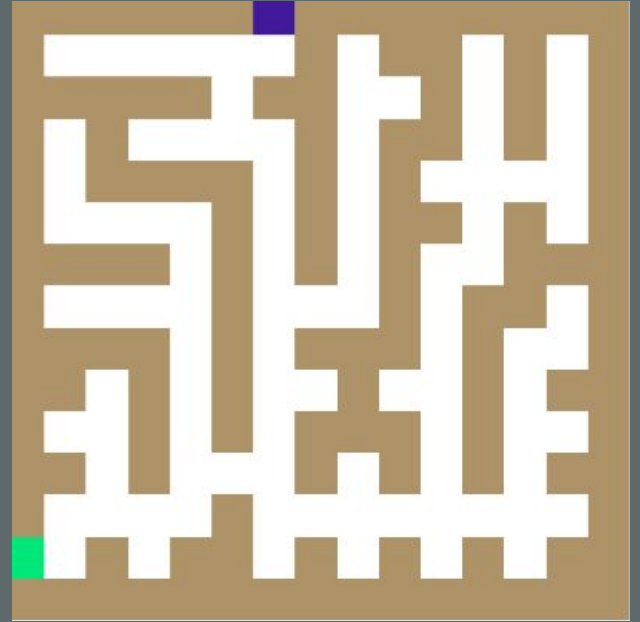


# Maze Solver with Reinforcement Learning

**Grigor Bezirganyan**  
**Henrik Sergoyan**

What is the project about?



# Steps

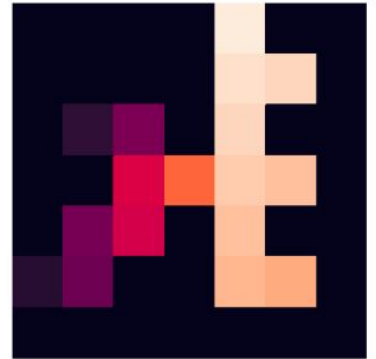
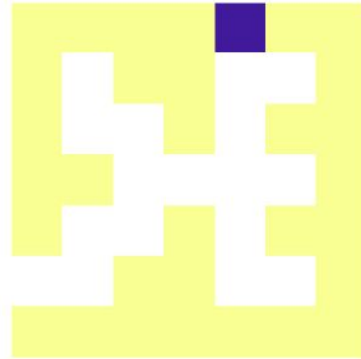
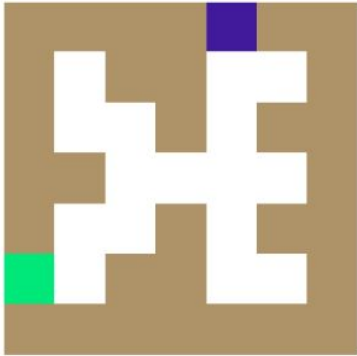
Background information

1. Creating the Environment
2. Learning Process
3. Developing the Optimal Policy
4. Solve Theeeeeeeee Maaaaaze!

---

# Creating the Environment

Epoch 5/5



# Learning Process

- Markov Property

$$p(\mathbf{s}_{t+1} | \mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_t, \mathbf{a}_t) = p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

- Reinforcement Learning

$$\mathbb{E}_{\tau_{\pi}} \left[ \sum_{t=0}^T r(\mathbf{s}_t, \mathbf{a}_t) \right] \rightarrow \max_{\pi}$$

$$\mathbb{E}_{\tau_{\pi}} \left[ \sum_{t=0}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right] \rightarrow \max_{\pi}, \quad 0 \leq \gamma \leq 1$$

discount factor

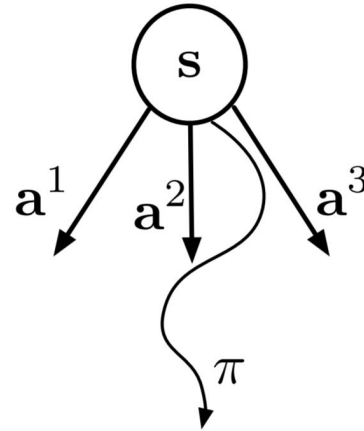
# Value Function

- Value function shows how good it is to act in accordance to policy  $\pi$  starting from some state  $s_t$

$$V^\pi(s_t) = \mathbb{E}_{\mathbf{a}_t, \mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \dots} [r(\mathbf{s}_t, \mathbf{a}_t) + \gamma r(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) + \dots]$$

- Optimal value function shows **the maximum amount of reward** we can get starting from some state  $s_t$

$$V^*(s_t) = \max_{\pi} V^\pi(s_t), \quad \forall s_t \in \mathcal{S}$$



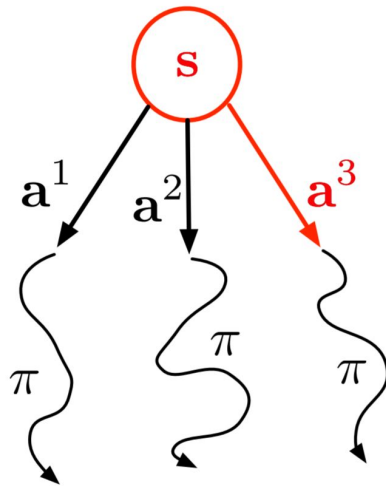
# Q-function

- Q-function shows how good it is to act in accordance to policy  $\pi$  starting from some state  $s_t$  after taking action  $a_t$

$$Q^\pi(s_t, a_t) = \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} [r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \dots]$$

- Optimal Q-function shows **the maximum amount of reward** we can get starting from some state  $s_t$  after taking action  $a_t$

$$Q^*(s_t, a_t) = \max_{\pi} Q^\pi(s_t, a_t), \quad \forall (s_t, a_t)$$



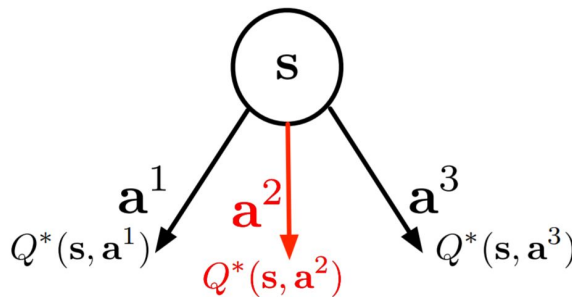
# Optimal Policy

Suppose we know the optimal Q-function.

$$Q^*(\mathbf{s}_t, \mathbf{a}_t) = \max_{\pi} Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t), \quad \forall(\mathbf{s}_t, \mathbf{a}_t)$$

What will be the optimal policy?

$$\pi^*(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1, & \mathbf{a}_t = \arg \max_{\mathbf{a}_t \in \mathcal{A}} Q^*(\mathbf{s}_t, \mathbf{a}_t) \\ 0, & \text{otherwise} \end{cases}$$





Thank You