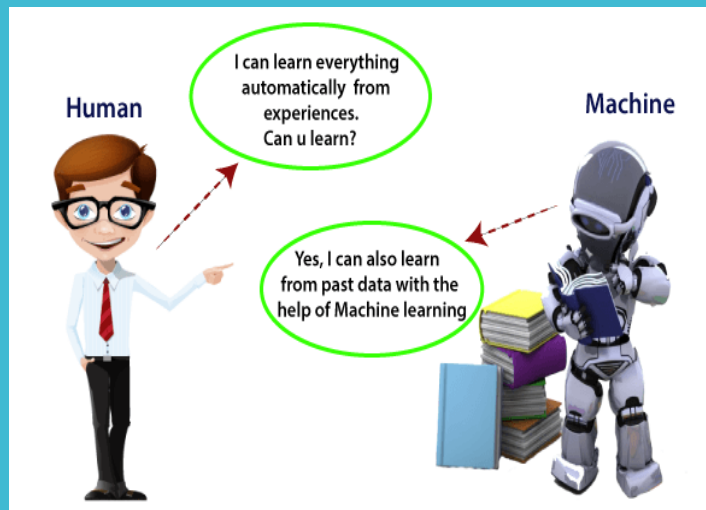
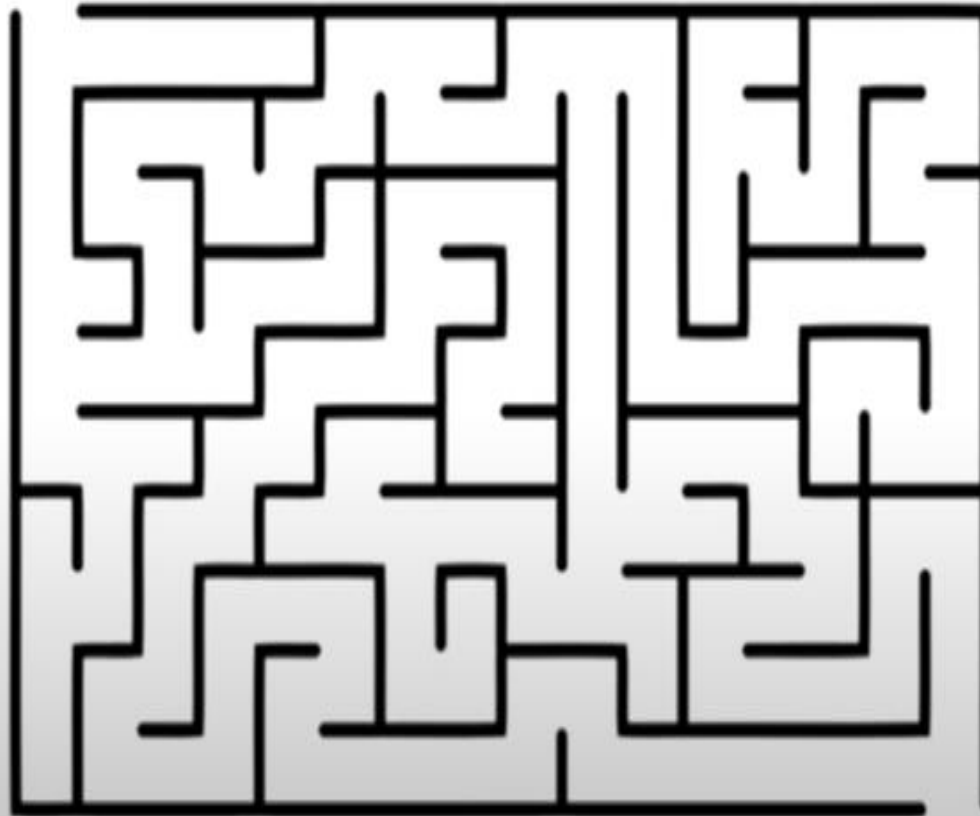
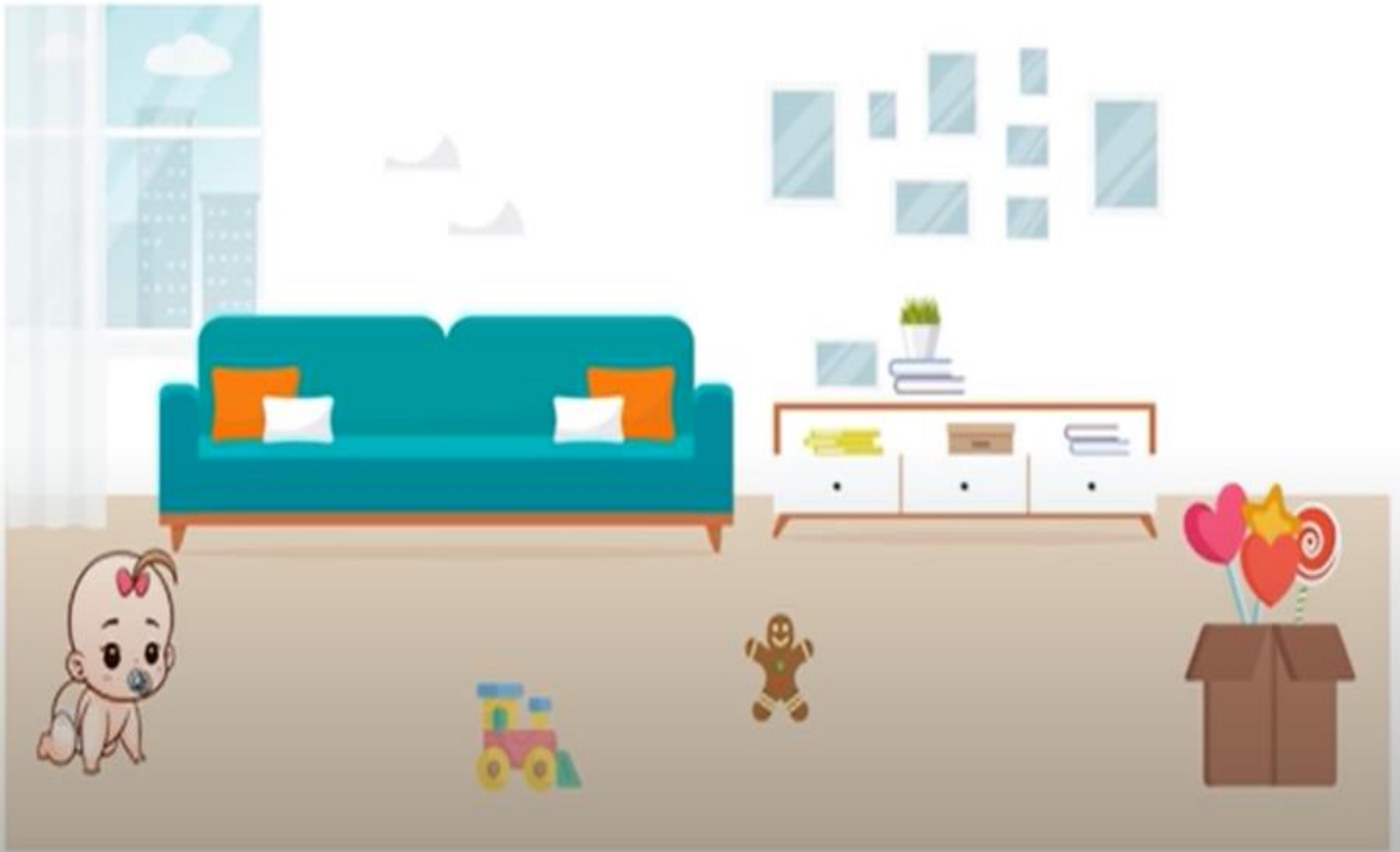


Reinforcement Learning

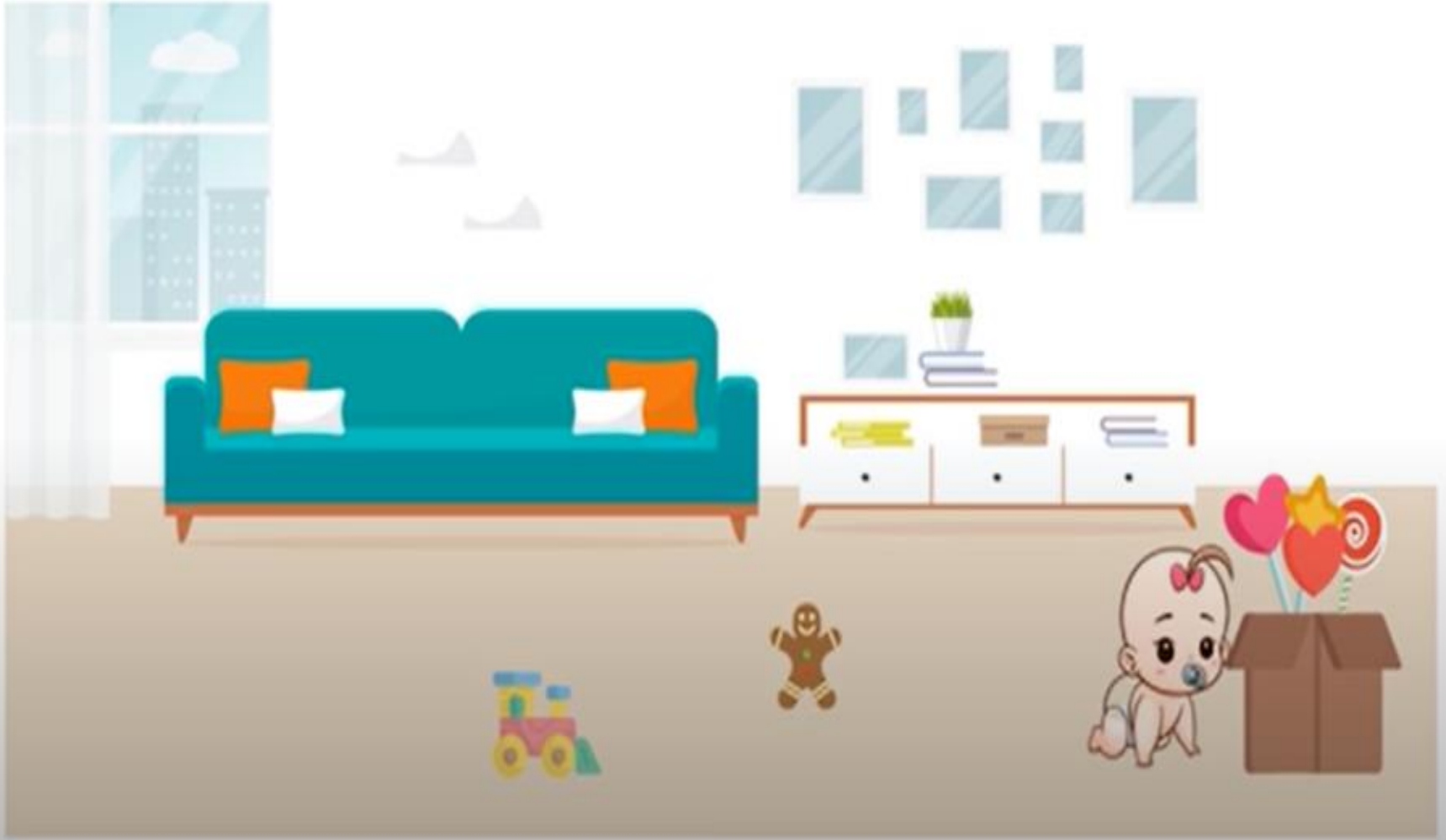


Reinforcement learning is a type of Machine Learning where an agent learns to behave in a environment by performing actions and seeing the results

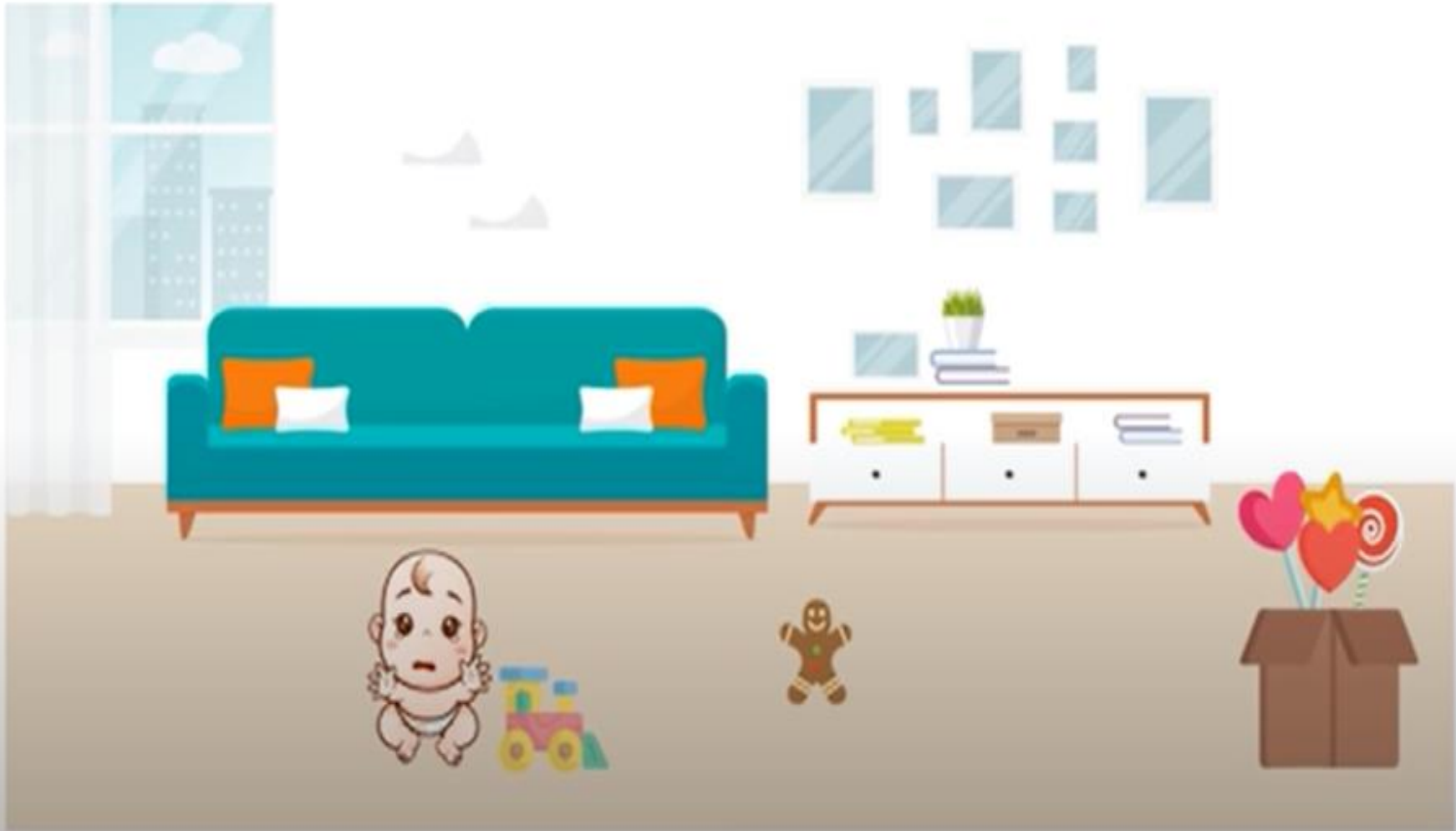




Scenario 1: Baby starts crawling and makes it to the candy



Scenario 2: Baby starts crawling but falls due to some hurdle in between



Reinforcement Learning system is comprised of two main components:

- Agent
- Environment

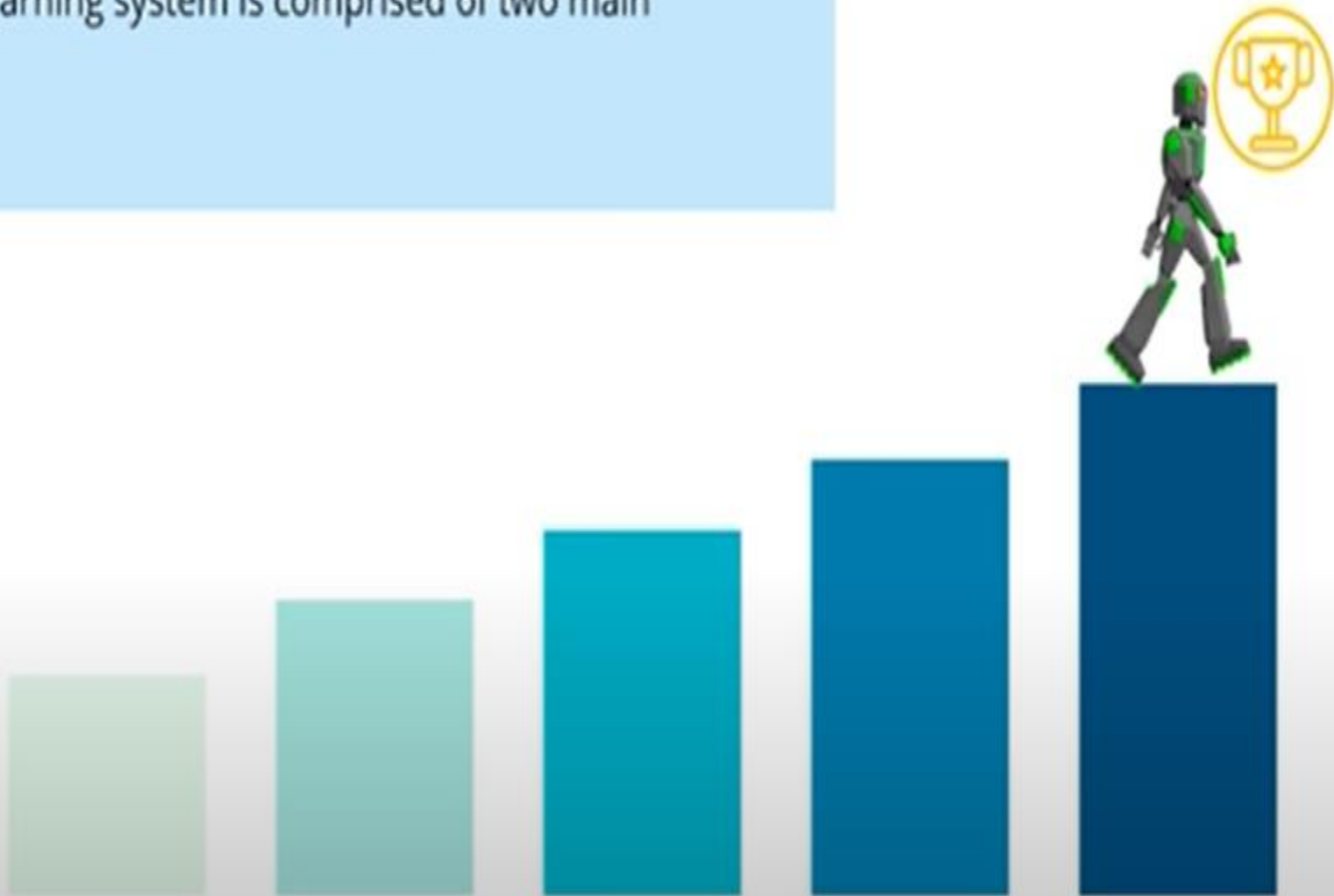
Agent



Environment

Reinforcement Learning system is comprised of two main components:

- Agent
- Environment





Agent: The RL algorithm that learns from trial and error

Environment: The world through which the agent moves



Action (A): All the possible steps that the agent can take

State (S): Current condition returned by the environment



Reward (R): An instant return from the environment to appraise the last action



Policy (π): The approach that the agent uses to determine the next action based on the current state

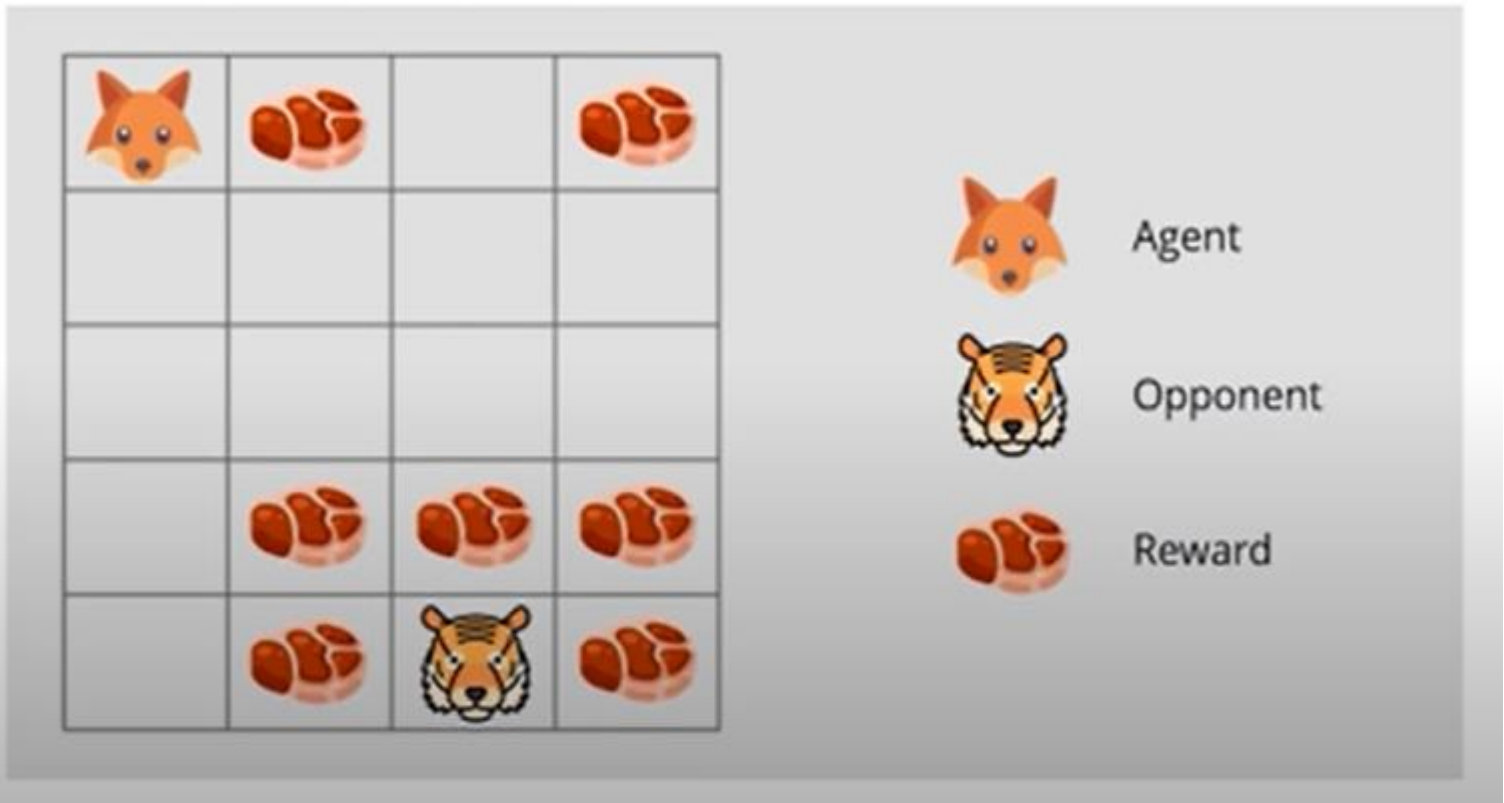


Value (V): The expected long-term return with discount, as opposed to the short-term reward R



Action-value (Q): This similar to Value, except, it takes an extra parameter, the current action (A)

Reward maximization theory states that, *a RL agent must be trained in such a way that, he takes the best action so that the reward is maximum.*



Exploration Vs Exploitation

Counter Strike Example

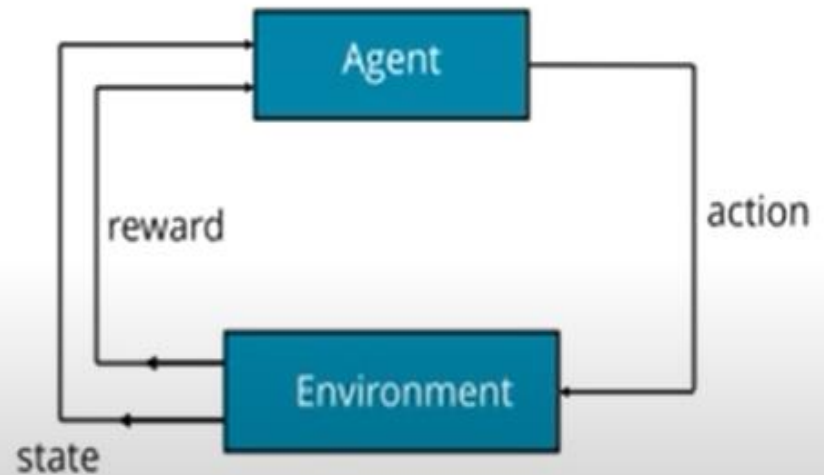


1. The RL Agent (Player1) collects state S^0 from the environment
2. Based on the state S^0 , the RL agent takes an action A^0 , initially the action is random
3. The environment is now in a new state S^1
4. RL agent now gets a reward R^1 from the environment
5. The RL loop goes on until the RL agent is dead or reaches the destination

The mathematical approach for mapping a solution in reinforcement learning is called *Markov Decision Process* (MDP)

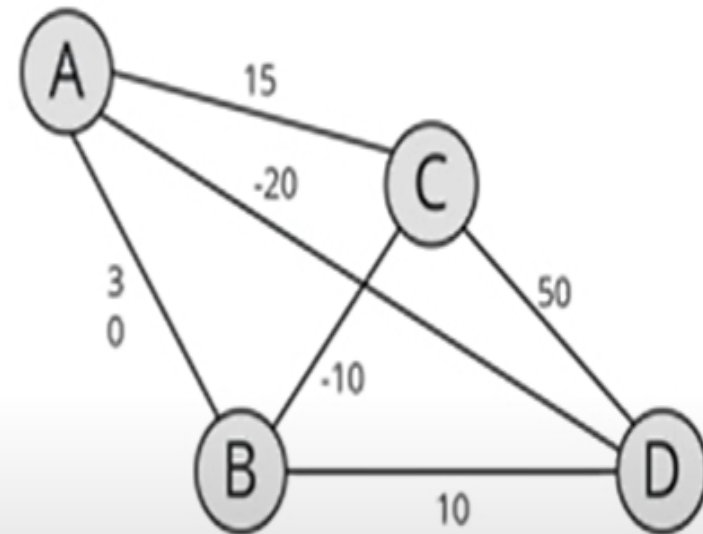
The following parameters are used to attain a solution:

- Set of actions, A
- Set of states, S
- Reward, R
- Policy, π
- Value, V

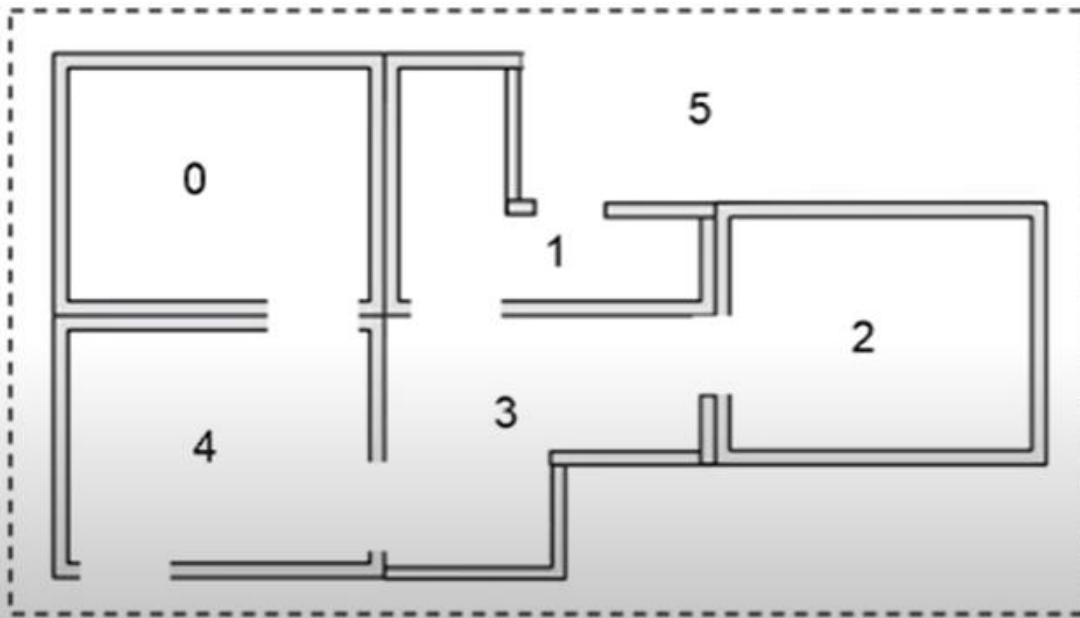


In this problem,

- Set of states are denoted by nodes i.e. {A, B, C, D}
- Action is to traverse from one node to another {A -> B, C -> D}
- Reward is the cost represented by each edge
- Policy is the path taken to reach the destination {A -> C -> D}



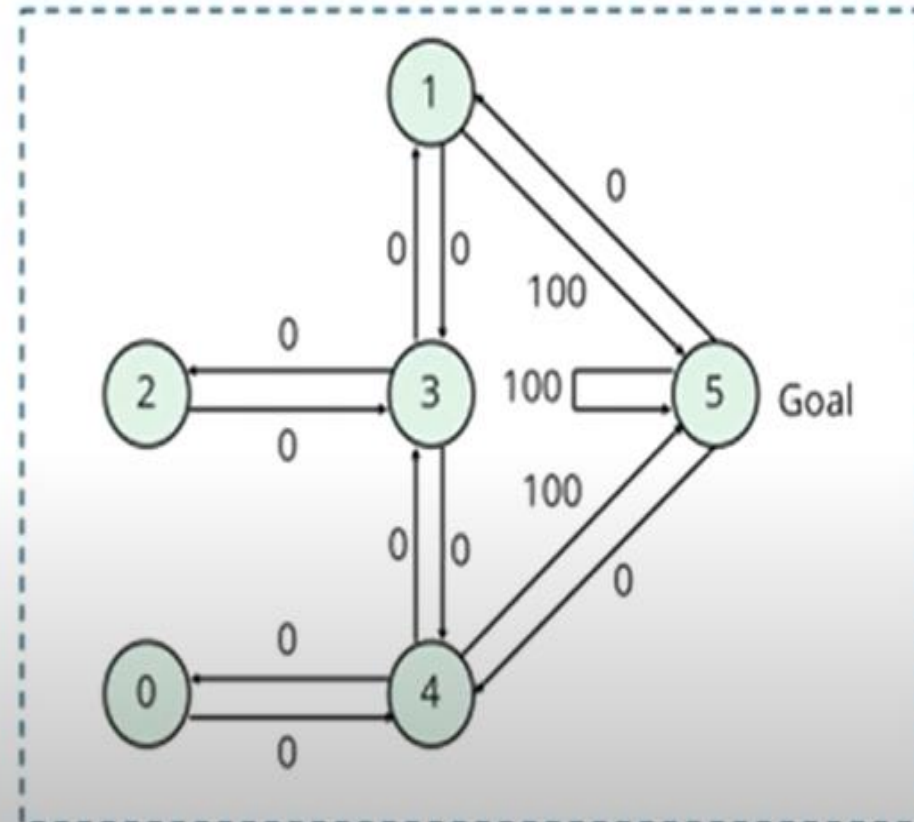
Place an agent in any one of the rooms (0,1,2,3,4) and the goal is to reach outside the building (room 5)



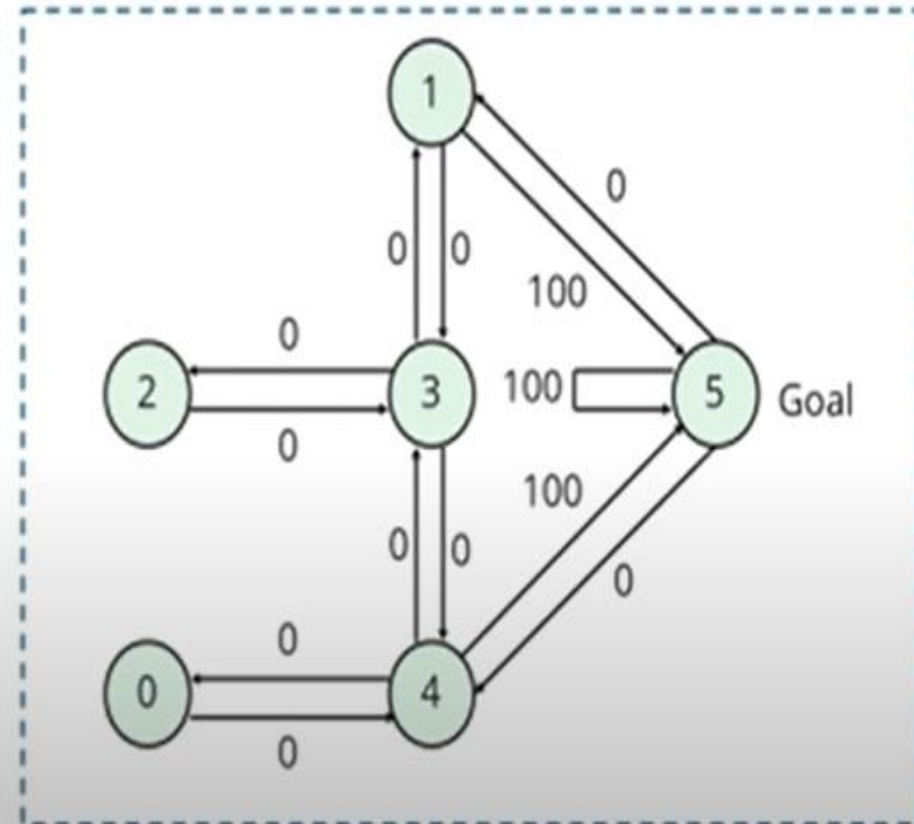
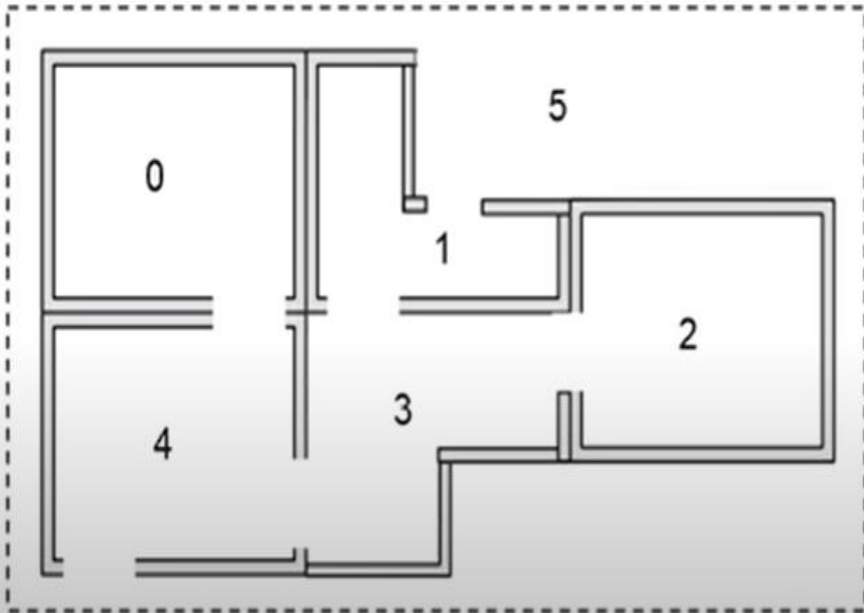
- 5 rooms in a building connected by doors
- each room is numbered 0 through 4
- The outside of the building can be thought of as one big room (5)
- Doors 1 and 4 lead into the building from room 5 (outside)

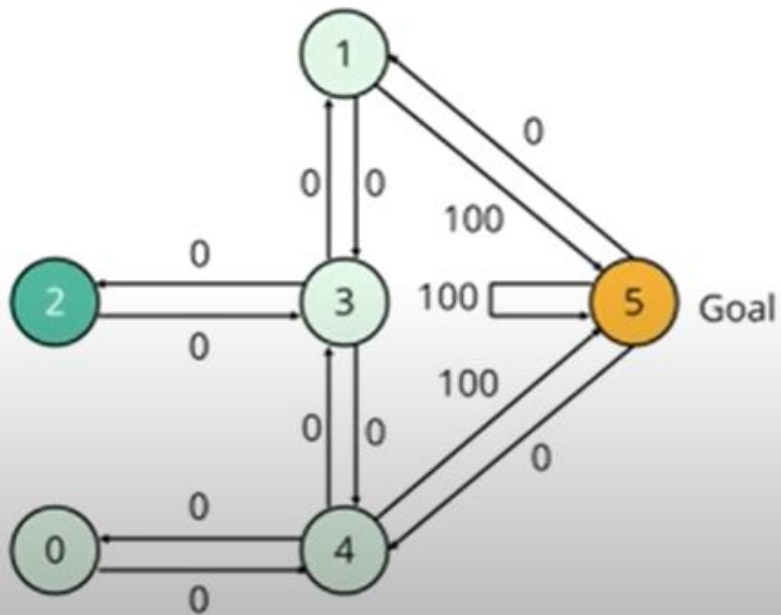
Next step is to associate a reward value to each door:

- doors that lead directly to the goal have a reward of 100
- Doors not directly connected to the target room have zero reward
- Because doors are two-way, two arrows are assigned to each room
- Each arrow contains an instant reward value



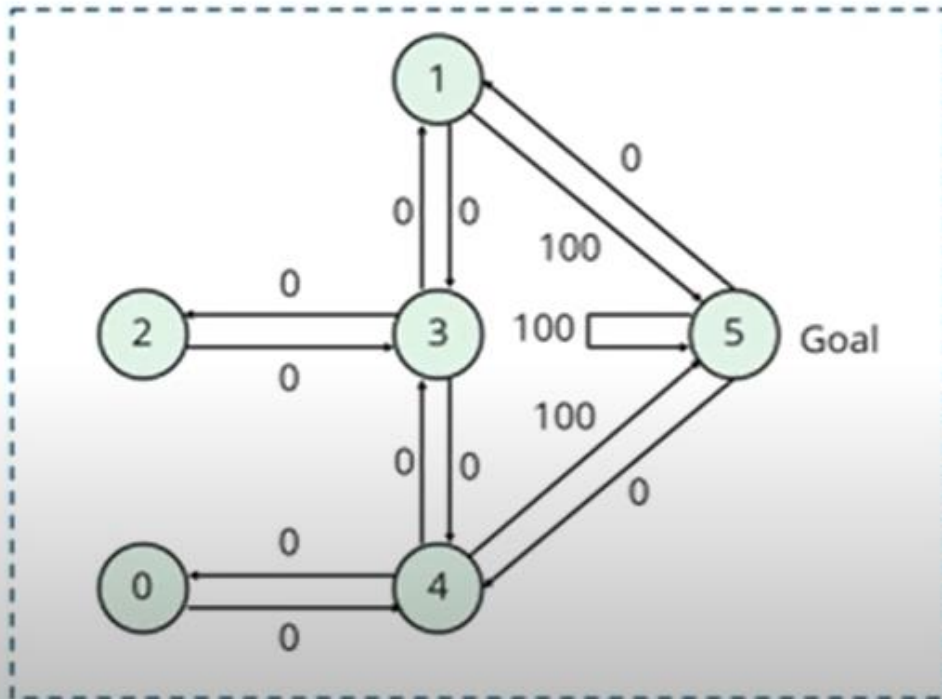
Place an agent in any one of the rooms (0,1,2,3,4) and the goal is to





Example (Agent traverse from room 2 to room5):

1. Initial state = state 2
2. State 2 \rightarrow state 3
3. State 3 \rightarrow state (2, 1, 4)
4. State 4 \rightarrow state 5



State

Action

0 1 2 3 4 5

0

1

2

3

4

5

$R =$

0	-1	-1	-1	-1	0	-1
1	-1	-1	-1	0	-1	100
2	-1	-1	-1	0	-1	-1
3	-1	0	0	-1	0	-1
4	0	-1	-1	0	-1	100
5	-1	0	-1	-1	0	100

The -1's in the table represent null values

Add another matrix Q, representing the memory of what the agent has learned through experience.

- The rows of matrix Q represent the current state of the agent
- columns represent the possible actions leading to the next state
- Formula to calculate the Q matrix:

$$Q(\text{state}, \text{action}) = R(\text{state}, \text{action}) + \text{Gamma} * \text{Max} [Q(\text{next state}, \text{all actions})]$$

Note

The Gamma parameter has a range of 0 to 1 ($0 \leq \text{Gamma} < 1$).

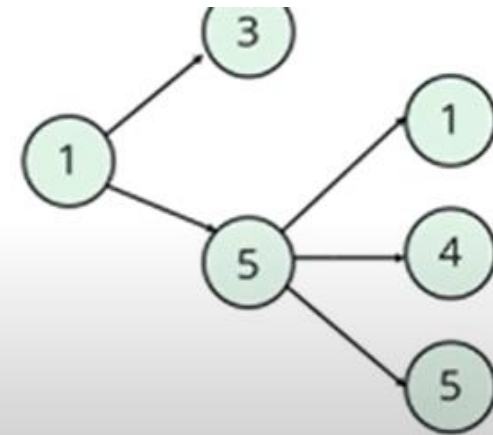
- If Gamma is closer to zero, the agent will tend to consider only immediate rewards.
- If Gamma is closer to one, the agent will consider future rewards with greater weight

- ① Set the gamma parameter, and environment rewards in matrix R
- ② Initialize matrix Q to zero
- ③ Select a random initial state
- ④ Set initial state = current state
- ⑤ Select one among all possible actions for the current state
- ⑥ Using this possible action, consider going to the next state
- ⑦ Get maximum Q value for this next state based on all possible actions
- ⑧ Compute: $Q(\text{state}, \text{action}) = R(\text{state}, \text{action}) + \text{Gamma} * \text{Max}[Q(\text{next state}, \text{all actions})]$
- ⑨ Repeat above steps until current state = goal state

$$Q(1,5) = R(1,5) + 0.8 * \text{Max}[Q(5,1), Q(5,4), Q(5,5)] = 100 + 0.8 * 0 = 100$$

		Action					
		0	1	2	3	4	5
Q =	0	0	0	0	0	0	0
	1	0	0	0	0	0	0
	2	0	0	0	0	0	0
	3	0	0	0	0	0	0
	4	0	0	0	0	0	0
	5	0	0	0	0	0	0

		Action					
		0	1	2	3	4	5
R =	0	-1	-1	-1	-1	0	-1
	1	-1	-1	-1	0	-1	100
	2	-1	-1	-1	0	-1	-1
	3	-1	0	0	-1	0	-1
	4	0	-1	-1	0	-1	100
	5	-1	0	-1	-1	0	100



$$Q(3,1) = R(3,1) + 0.8 * \text{Max}[Q(1,3), Q(1,5)] = 0 + 0.8 * [0, 100] = 80$$

The matrix Q get's updated

		Action					
		0	1	2	3	4	5
Q =	0	0	0	0	0	0	0
	1	0	0	0	0	0	100
	2	0	0	0	0	0	0
	3	0	80	0	0	0	0
	4	0	0	0	0	0	0
	5	0	0	0	0	0	0

		Action					
		0	1	2	3	4	5
R =	0	-1	-1	-1	-1	0	-1
	1	-1	-1	-1	0	-1	100
	2	-1	-1	-1	0	-1	-1
	3	-1	0	0	-1	0	-1
	4	0	-1	-1	0	-1	100
	5	-1	0	-1	-1	0	100

