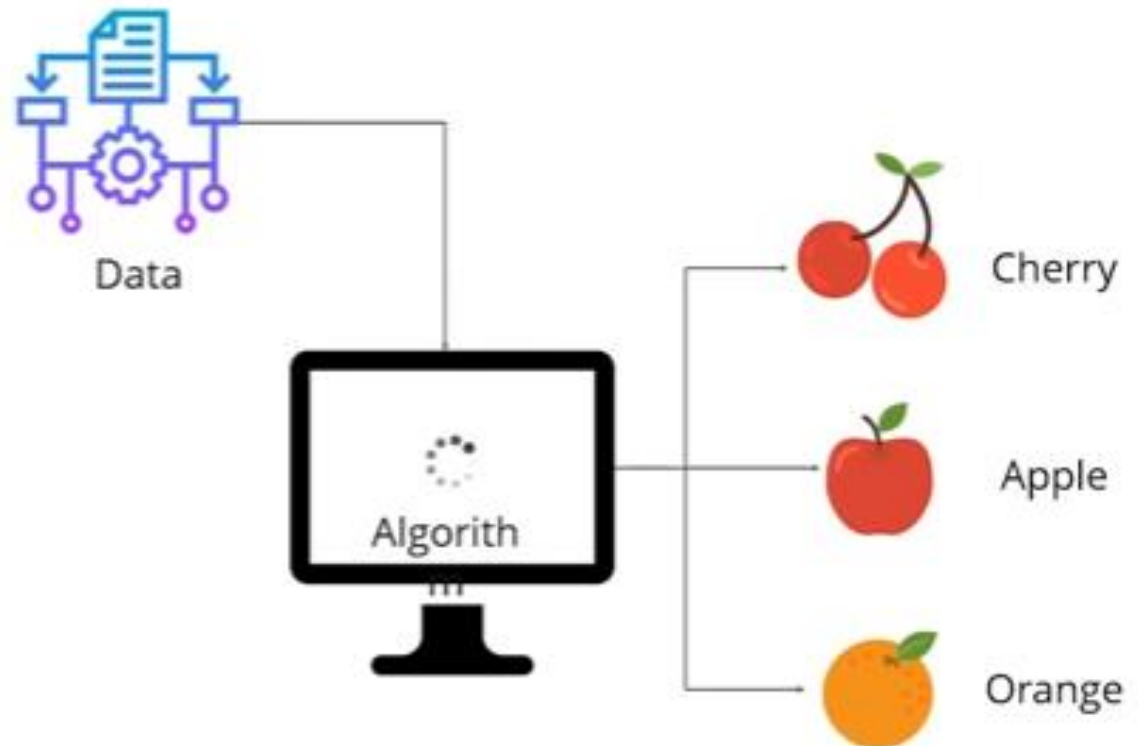# Reinforcement Learning

BDA/ HDA

# What Is Machine Learning?

Machine learning is a subset of artificial intelligence (AI) which provides machines the ability to learn automatically & improve from experience without being explicitly programmed.



They look the same!

Data

Algorith m

Cherry

Apple

Orange

# Types Of Machine Learning
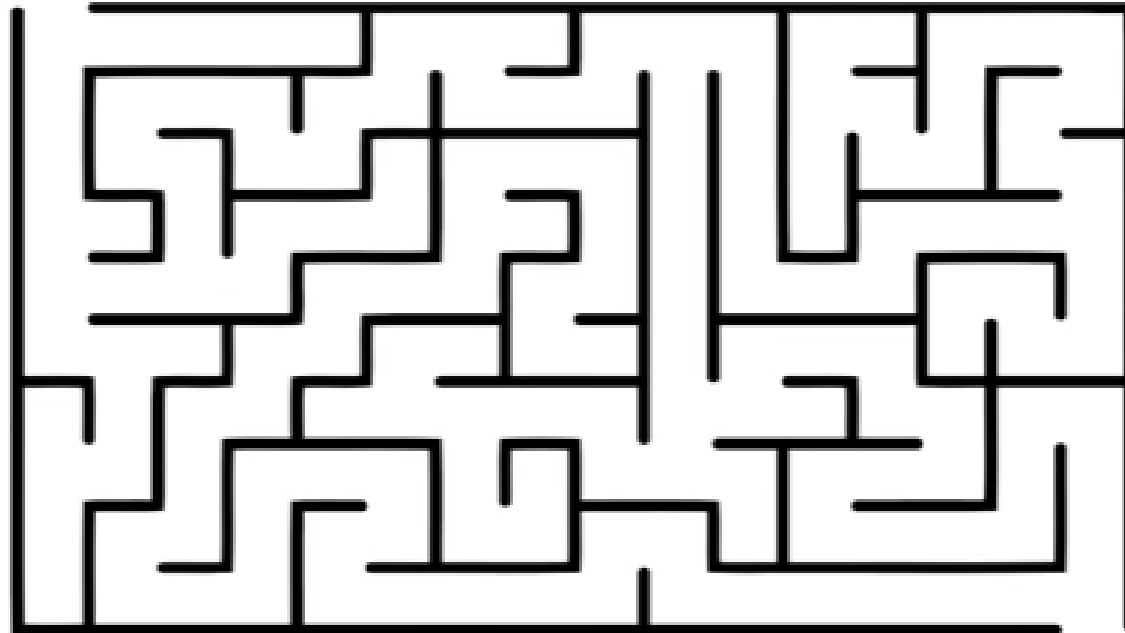


Supervised Learning

Unsupervised Learning

Reinforcement Learning

# What Is 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
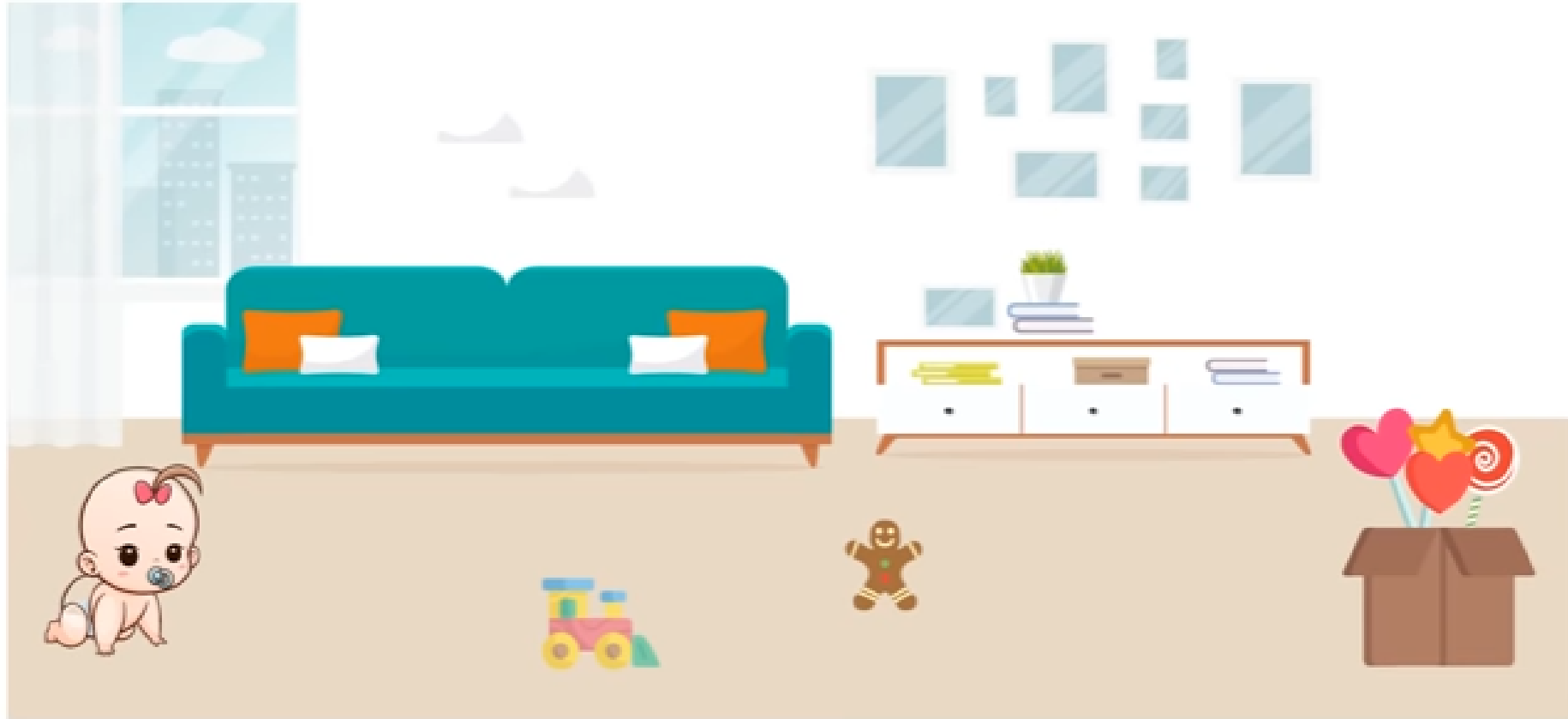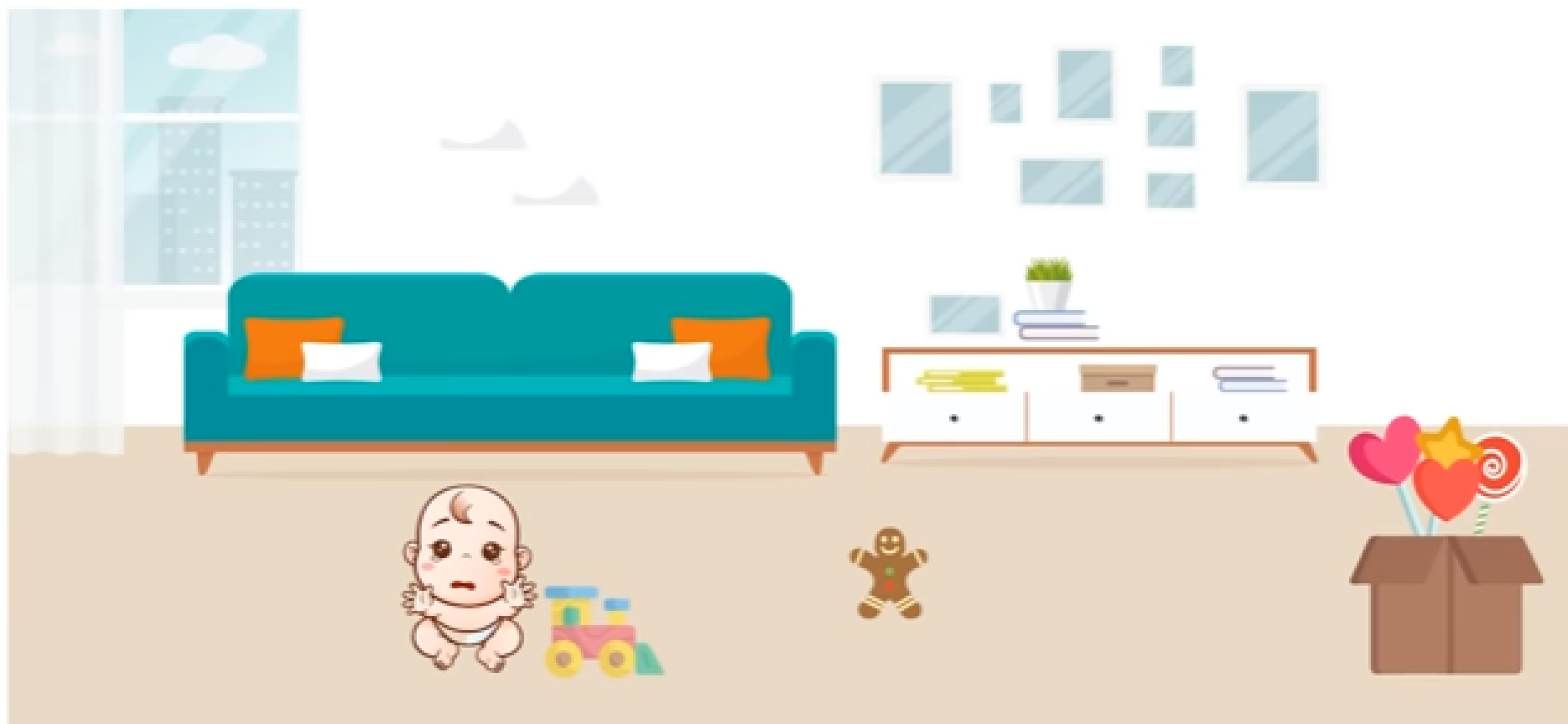
Agent

Environment

Reward

# Reinforcement Learning With An Analogy

# Reinforcement Learning With An Analogy
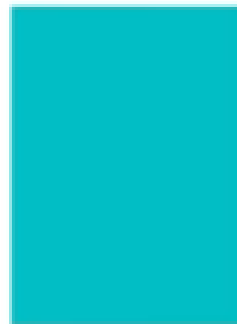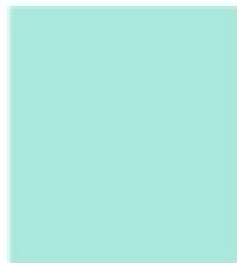
# Reinforcement Learning vs. Supervised Learning

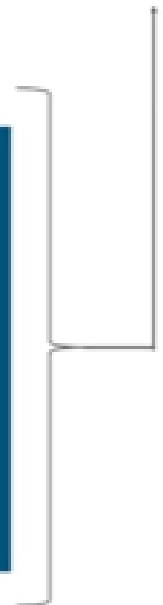| Parameters | Reinforcement Learning | Supervised Learning |
| --- | --- | --- |
| Decision style | Take decisions sequentially. | Decision is made on the input given at the beginning. |
| Works on | Works on interacting with the environment. | Works on examples or given sample data. |
| Dependency on decision | Decision is dependent. | Decisions are independent of each other. |
| Best suited | Supports and work better in AI, where human interaction is prevalent. | It is mostly operated with an interactive software system or applications. |
| Example | Chess game | Object recognition |

# Reinforcement Learning Process

Reinforcement Learning system is comprised of two main components:
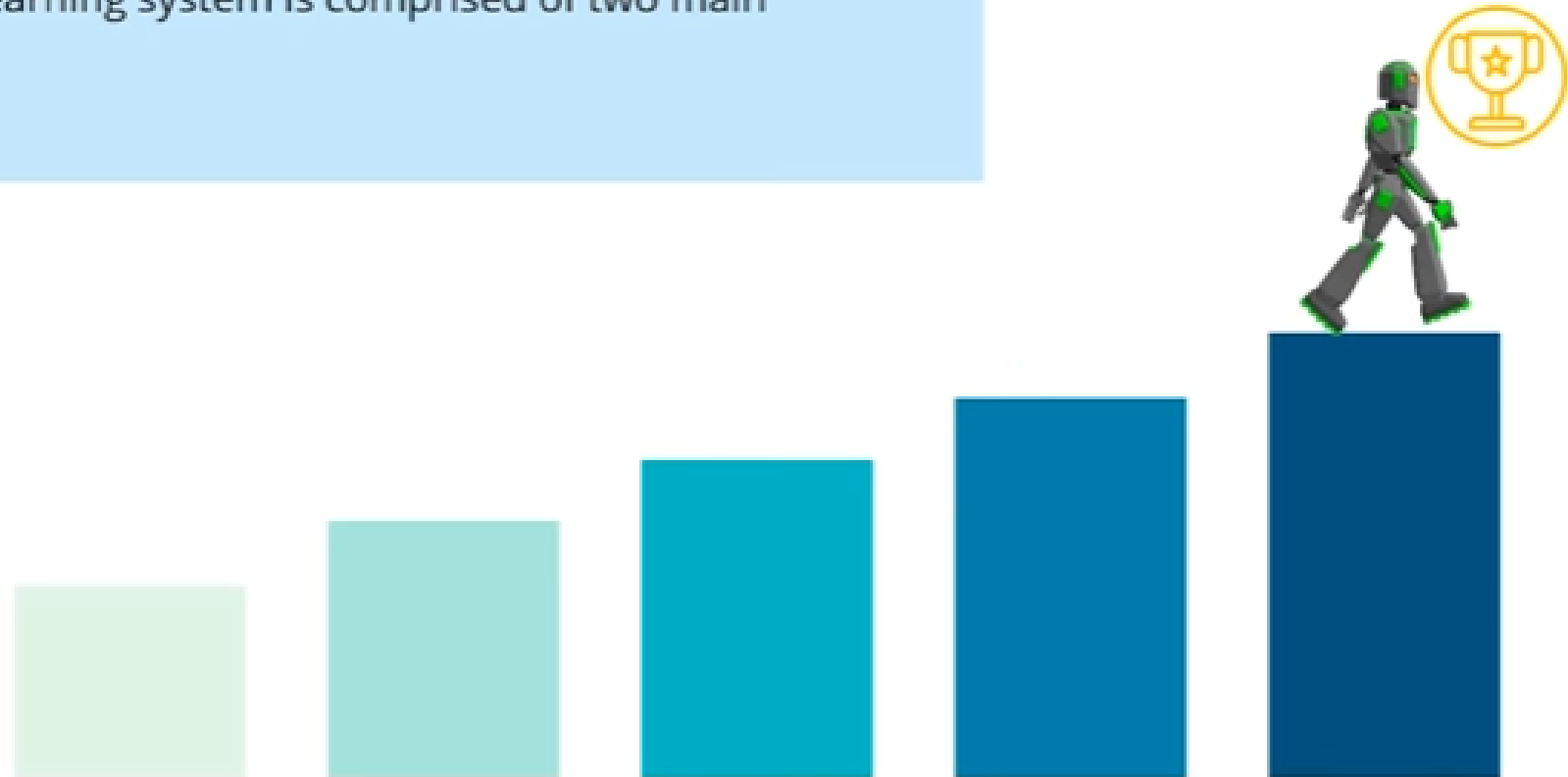- Agent
- Environment

Environment

Agent

# Reinforcement Learning Process

Reinforcement Learning system is comprised of two main components:
- Agent
- Environment

# 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

# Reinforcement Learning Definitions

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

# Reinforcement Learning Definitions

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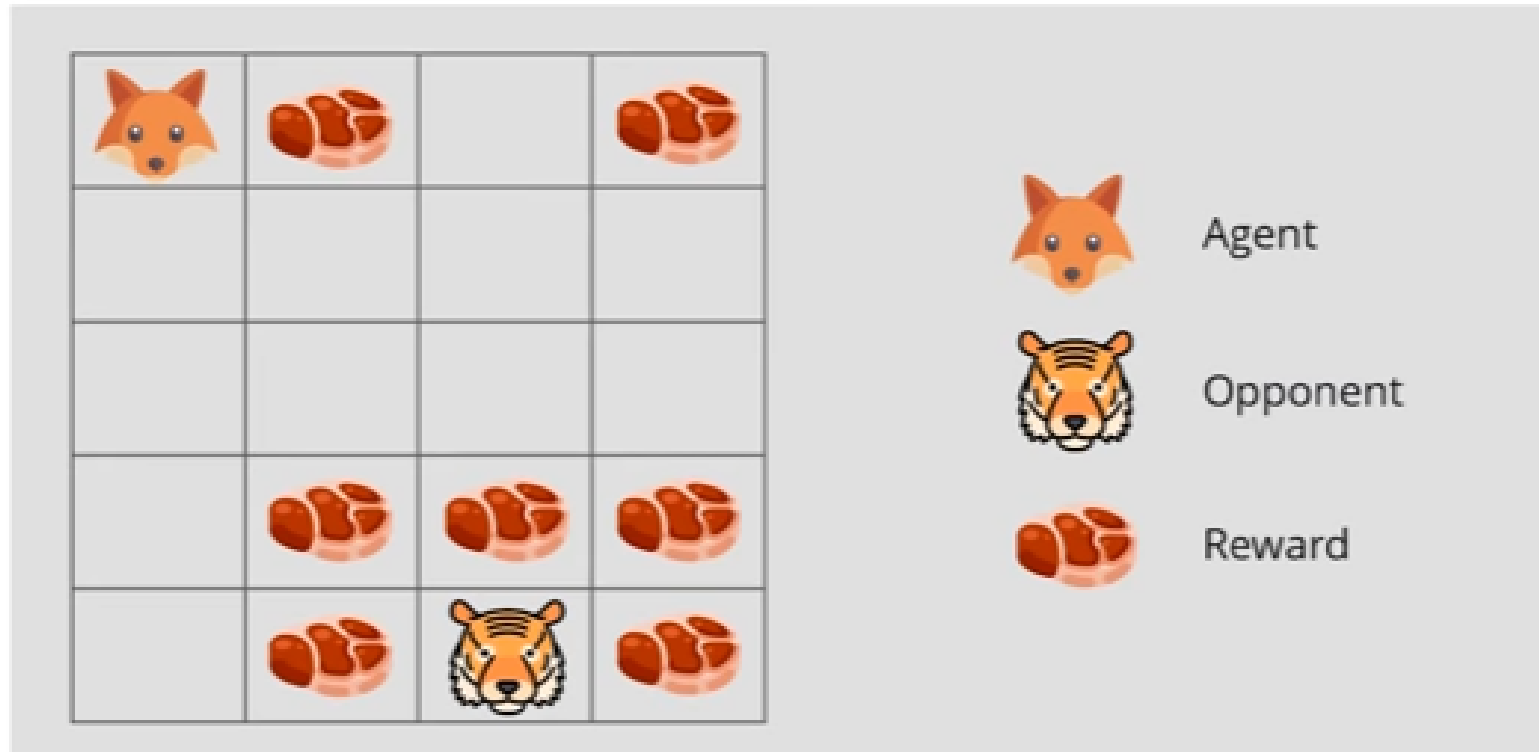
Value (V): The expected long-term return with discount, as opposed to the short-term reward R

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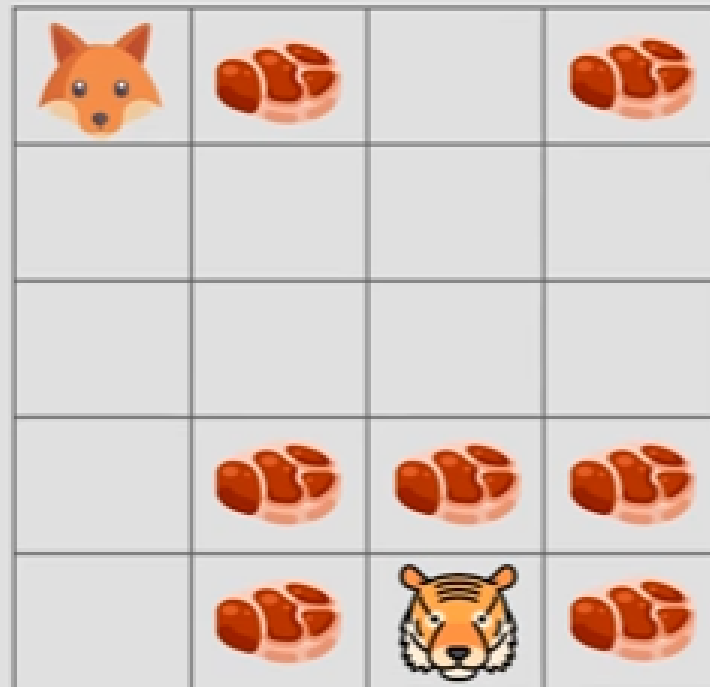
# Reward Maximization

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 & Exploitation

*Exploitation* is about using the already known exploited information to heighten the rewards

*Exploration* is about exploring and capturing more information about an environment
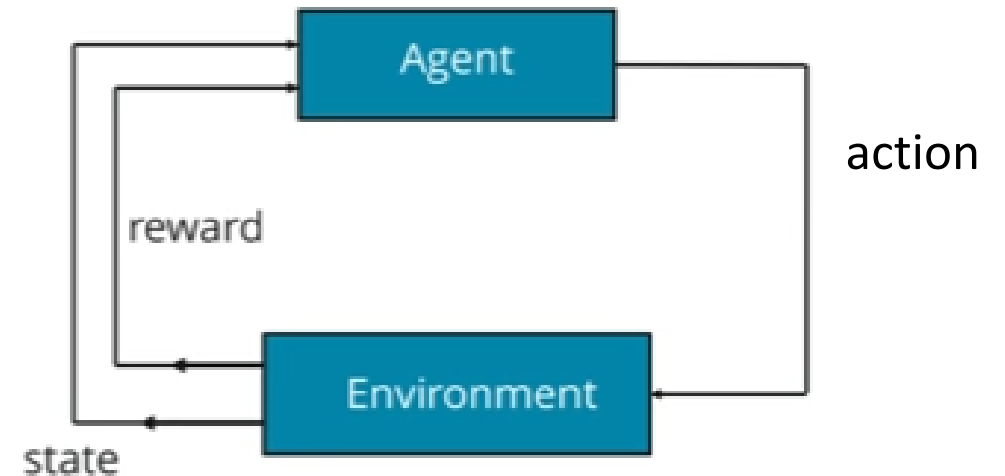
# Markov Decision Process

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
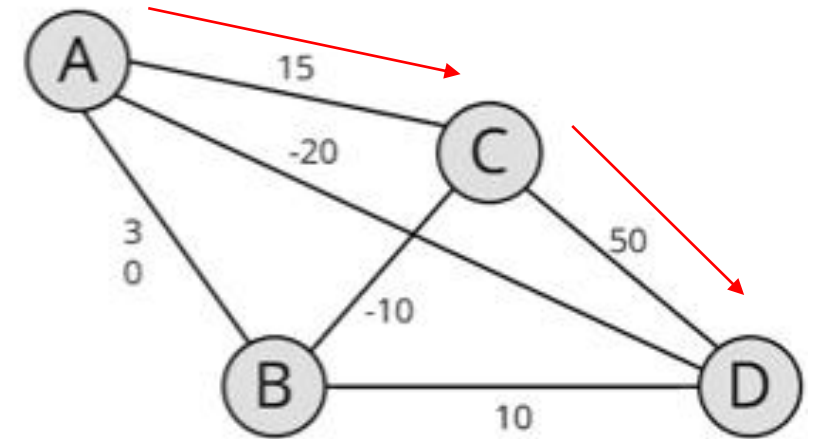
# Markov Decision Process – Shortest Path Problem

Goal: Find the shortest path between A and D with minimum possible cost (maximum reward)

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}

# Understanding Q-Learning With An Example

*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)

# Understanding Q-Learning With An Example

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

# Understanding Q-Learning With An Example

The terminology in Q-Learning includes the terms state and action:
- Room (including room 5) represents a state
- agent's movement from one room to another represents an action
- In the figure, a state is depicted as a node, while "action" is represented by the arrows



Example (Agent traverse from room 2 to room5):

1. Initial state = state 2

2. State 2 -> state 3

3. State 3 -> state (2, 1, 4)

4. State 4 -> state 5

# Understanding Q-Learning With An Example

We can put the state diagram and the instant reward values into a reward table, matrix R.



$$R = \begin{array}{c|cccccc} & 0 & 1 & 2 & 3 & 4 & 5 \\ \hline 0 & -1 & -1 & -1 & -1 & 0 & -1 \\ 1 & -1 & -1 & -1 & 0 & -1 & 100 \\ 2 & -1 & -1 & -1 & 0 & -1 & -1 \\ 3 & -1 & 0 & 0 & -1 & 0 & -1 \\ 4 & 0 & -1 & -1 & 0 & -1 & 100 \\ 5 & -1 & 0 & -1 & -1 & 0 & 100 \end{array}$$

The -1's in the table represent null values

# Understanding Q-Learning With An Example

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:

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

**Note**

The Gamma parameter has a range of 0 to 1 (0 <= 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

# Q – Learning Example

First step is to set the value of the learning parameter Gamma = 0.8, and the initial state as Room 1.
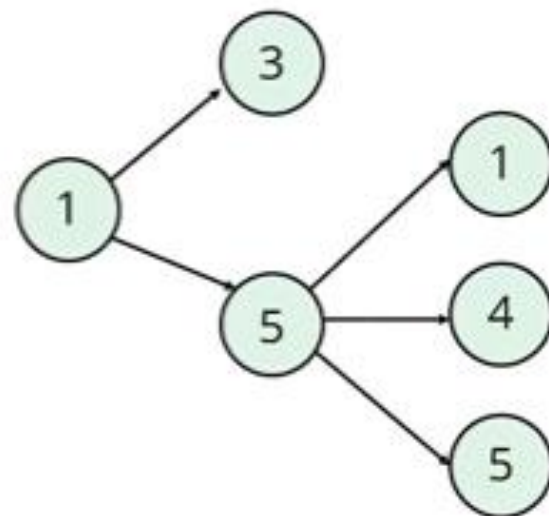
Next, initialize matrix Q as a zero matrix:
- From room 1 you can either go to room 3 or 5, let's select room 5.
- From room 5, calculate maximum Q value for this next state based on all possible actions:

Q(state, action) = R(state, action) + Gamma * Max[Q(next state, all actions)]

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

|     | 0 | 1 | 2 | 3 | 4 | 5 |
|-----|---|---|---|---|---|---|
| 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 |

Q =

Action

| State | 0 | 1 | 2 | 3 | 4 | 5 |
|-------|----|----|----|----|----|-----|
| 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 |

R =

# Q – Learning Algorithm

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(state, action) = R(state, action) + Gamma * Max[Q(next state, all actions)]

⑨ Repeat above steps until current state = goal state

# Q – Learning Example

First step is to set the value of the learning parameter Gamma = 0.8, and the initial state as Room 1.

Next, initialize matrix Q as a zero matrix:
- From room 1 you can either go to room 3 or 5, let's select room 5.
- From room 5, calculate maximum Q value for this next state based on all possible actions:

Q(state, action) = R(state, action) + Gamma * Max[Q(next state, all actions)]

Q(1,5) = R(1,5) + 0.8 * Max[Q(5,1), Q(5,4), Q(5,5)] = 100 + 0.8 * 0 = 100

$$Q = \begin{array}{c|cccccc} & 0 & 1 & 2 & 3 & 4 & 5 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 100 \\ 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 3 & 0 & 0 & 0 & 0 & 0 & 0 \\ 4 & 0 & 0 & 0 & 0 & 0 & 0 \\ 5 & 0 & 0 & 0 & 0 & 0 & 0 \end{array}$$

Action

$$R = \begin{array}{c|cccccc} \text{State} & 0 & 1 & 2 & 3 & 4 & 5 \\ \hline 0 & -1 & -1 & -1 & -1 & 0 & -1 \\ 1 & -1 & -1 & -1 & 0 & -1 & 100 \\ 2 & -1 & -1 & -1 & 0 & -1 & -1 \\ 3 & -1 & 0 & 0 & -1 & 0 & -1 \\ 4 & 0 & -1 & -1 & 0 & -1 & 100 \\ 5 & -1 & 0 & -1 & -1 & 0 & 100 \end{array}$$

# Q – Learning Example

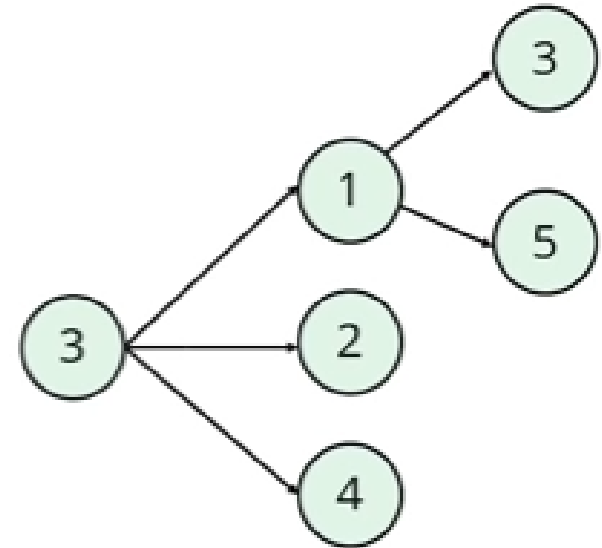For the next episode, we start with a randomly chosen initial state, i.e. state 3
- From room 3 you can either go to room 1,2 or 4, let's select room 1.
- From room 1, calculate maximum Q value for this next state based on all possible actions:

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

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

The matrix Q get's updated

$$Q = \begin{array}{c|cccccc}
 & 0 & 1 & 2 & 3 & 4 & 5 \\
\hline
0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 0 & 0 & 0 & 0 & 0 & 100 \\
2 & 0 & 0 & 0 & 0 & 0 & 0 \\
3 & 0 & 80 & 0 & 0 & 0 & 0 \\
4 & 0 & 0 & 0 & 0 & 0 & 0 \\
5 & 0 & 0 & 0 & 0 & 0 & 0 \\
\end{array}$$

$$R = \begin{array}{c|cccccc}
\text{State} & 0 & 1 & 2 & 3 & 4 & 5 \\
\hline
0 & -1 & -1 & -1 & -1 & 0 & -1 \\
1 & -1 & -1 & -1 & 0 & -1 & 100 \\
2 & -1 & -1 & -1 & 0 & -1 & -1 \\
3 & -1 & 0 & 0 & -1 & 0 & -1 \\
4 & 0 & -1 & -1 & 0 & -1 & 100 \\
5 & -1 & 0 & -1 & -1 & 0 & 100 \\
\end{array}$$

Action

# Thank you