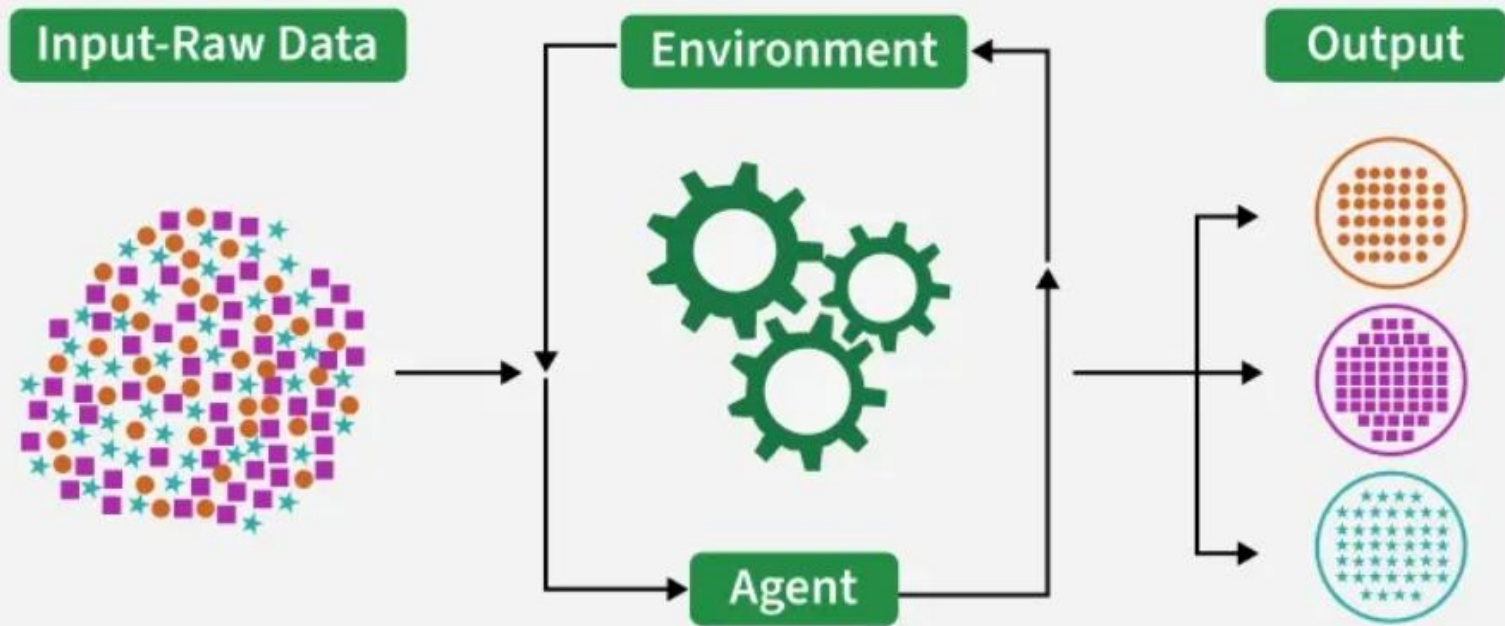


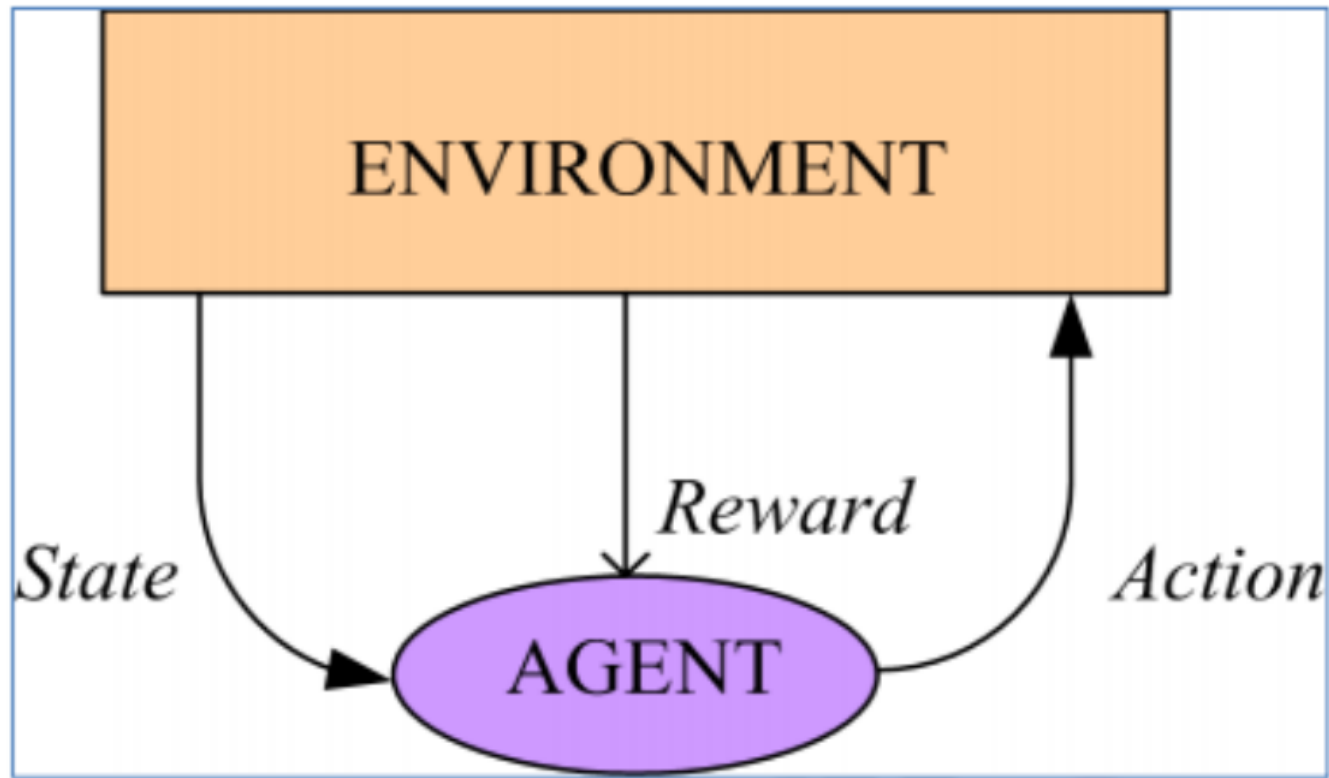
# Reinforcement Learning

## Reinforcement Learning in ML



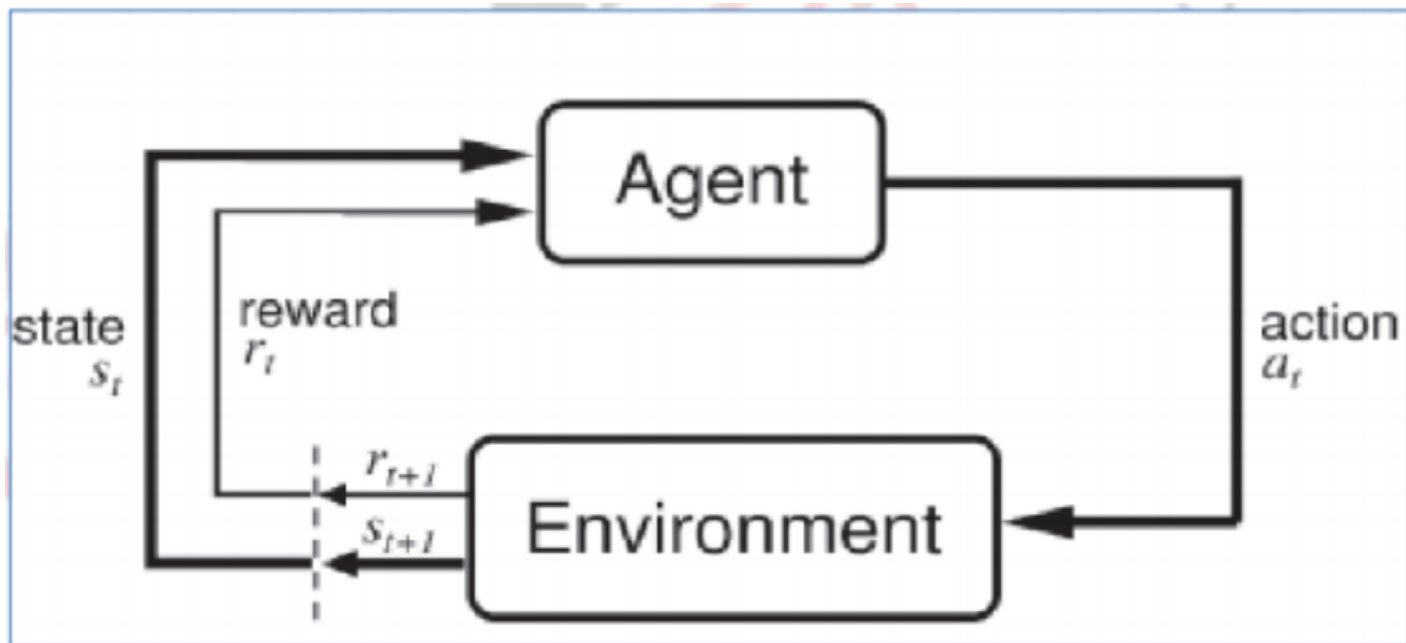
# Reinforcement Learning

- ❑ Reinforcement Learning (RL) is a branch of machine learning that focuses on how agents can learn to make decisions through trial and error to maximize cumulative rewards.
- ❑ RL allows machines to learn by interacting with an environment and receiving feedback based on their actions. This feedback comes in the form of rewards or penalties.



The core of Reinforcement Learning includes contributors such as:

- ❑ **Agent:** The decision-maker that performs actions.
- ❑ **Environment:** The world or system in which the agent operates.
- ❑ **State:** The situation or condition the agent is currently in.
- ❑ **Action:** The possible moves or decisions the agent can make.
- ❑ **Reward:** The feedback or result from the environment based on the agent's action.



# Core Components

Let's see the core components of Reinforcement Learning

## 1. Policy

- ❑ Defines the agent's behavior i.e maps states for actions.
- ❑ Can be simple rules or complex computations.

**Example:** An autonomous car maps pedestrian detection to make necessary stops.

## 2. Reward Signal

- ❑ Represents the goal of the RL problem.
- ❑ Guides the agent by providing feedback (positive/negative rewards).

**Example:** For self-driving cars rewards can be fewer collisions, shorter travel time, lane discipline.





### 3. Value Function

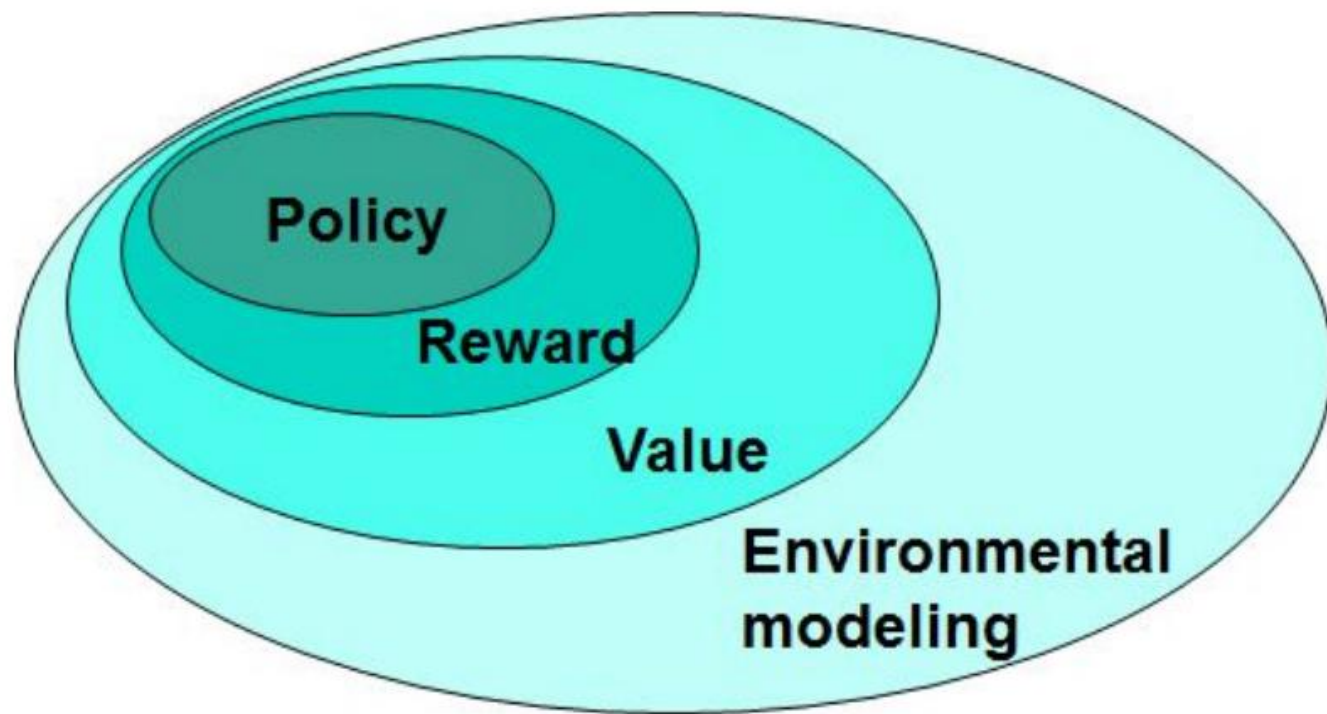
- ❑ Evaluates long-term benefits, not just immediate rewards.
- ❑ Measures desirability of a state considering future outcomes.

**Example:** A vehicle may avoid reckless maneuvers (short-term gain) to maximize overall safety and efficiency.

## **4. Model**

- ❑ Simulates the environment to predict outcomes of actions.
- ❑ Enables planning and foresight.

**Example:** Predicting other vehicles' movements to plan safer routes.



# Working of Reinforcement Learning

The agent interacts iteratively with its environment in a feedback loop:

- ❑ The agent observes the current state of the environment.
- ❑ It chooses and performs an action based on its policy.
- ❑ The environment responds by transitioning to a new state and providing a reward (or penalty).

- ❑ The agent updates its knowledge (policy, value function) based on the reward received and the new state.
- ❑ This cycle repeats with the agent balancing exploration (trying new actions) and exploitation (using known good actions) to maximize the cumulative reward over time.

# Types of Reinforcements

- ❑ Positive Reinforcement: Positive Reinforcement is defined as when an event, occurs due to a particular behavior, increases the strength and the frequency of the behavior. In other words, it has a positive effect on behavior.
- Advantages: Maximizes performance, helps sustain change over time.
- Disadvantages: Overuse can lead to excess states that may reduce effectiveness.

- ❑ Negative Reinforcement: Negative Reinforcement is defined as strengthening of behavior because a negative condition is stopped or avoided.
- Advantages: Increases behavior frequency, ensures a minimum performance standard.
- Disadvantages: It may only encourage just enough action to avoid penalties.



# Types of Reinforcement Learning

*Value-Based  
Reinforcement  
Learning*

*Policy-Based  
Reinforcement  
Learning*

*Model-Based  
Reinforcement  
Learning*

*Hybrid  
Approaches:  
Actor-Critic  
Methods*

# Value-Based Reinforcement Learning

- ❑ Value-based reinforcement learning focuses on finding the optimal value function that measures how good it is for an agent to be in a given state (or take a given action).
- ❑ The goal is to maximize the value function, which represents the long-term cumulative reward.
- ❑ The most common technique in this category is Q-learning.

# Policy-Based Reinforcement Learning

- ❑ Unlike value-based methods, policy-based RL methods aim to directly learn the optimal policy  $\pi(a|s)$ , which maps states to probabilities of selecting actions.
- ❑ These methods can be effective for environments with high-dimensional or continuous action spaces, where value-based methods struggle.

## ➤ **Advantages of Policy-Based Methods:**

- ❑ Effective in high-dimensional or continuous action spaces.
- ❑ Can learn stochastic policies, which can be beneficial in environments requiring exploration.

## ➤ **Challenges of Policy-Based Methods:**

- ❑ High variance in the gradient estimates.
- ❑ Often requires careful tuning of learning rates and other hyperparameters.

# Proximal Policy Optimization

- ❑ PPO is an improvement over basic policy gradient methods.
- ❑ It introduces a more stable way of updating the policy by clipping the update to prevent drastic changes.
- ❑ This ensures a more robust learning process, making PPO one of the most widely used algorithms in policy-based reinforcement learning.

# Model-Based Reinforcement Learning

- ❑ Model-based RL introduces an explicit model of the environment to predict the future states and rewards.
- ❑ The agent uses the model to simulate different actions and their outcomes before actually interacting with the environment.
- ❑ This helps the agent plan actions more effectively.

# Markov Decision Process

- ❑ This process is mathematically framed as a Markov Decision Process (MDP) where future states depend only on the current state and action, not on the prior sequence of events.
- ❑ Markov Decision Process (MDP) is a way to describe how a decision-making agent like a robot or game character moves through different situations while trying to achieve a goal.

❑ MDPs rely on variables such as the environment, agent's actions and rewards to decide the system's next optimal action. It helps us answer questions like:

- What actions should the agent take?
- What happens after an action?
- Is the result good or bad?



- ❑ In artificial intelligence Markov Decision Processes (MDPs) are used to model situations where decisions are made one after another and the results of actions are uncertain.
- ❑ They help in designing smart machines or agents that need to work in environments where each action might lead to different outcomes.

# Key Components of an MDP

- ❑ An MDP has five main parts:

## Markov Decision Process

States:  $S$

Model:  $T(S, a, S') \sim P(S' | S, a)$

Actions:  $A(S), A$

Reward:  $R(S), R(S, a), R(S, a, S')$

---

Policy:  $\Pi(S) \rightarrow a$   
 $\Pi^\star$

*Components of Markov Decision Process*

❑ **States (S):** A state is a situation or condition the agent can be in.

For example, A position on a grid like being at cell (1,1).

❑ **Actions (A):** An action is something the agent can do. For example, Move UP, DOWN, LEFT or RIGHT. Each state can have one or more possible actions.

❑ **Transition Model (T):** The model tells us what happens when an action is taken in a state. It's like asking: "If I move RIGHT from here, where will I land?" Sometimes the outcome isn't always the same that's uncertainty. For example:

- 80% chance of moving in the intended direction
- 10% chance of slipping to the left
- 10% chance of slipping to the right

This randomness is called a stochastic transition.

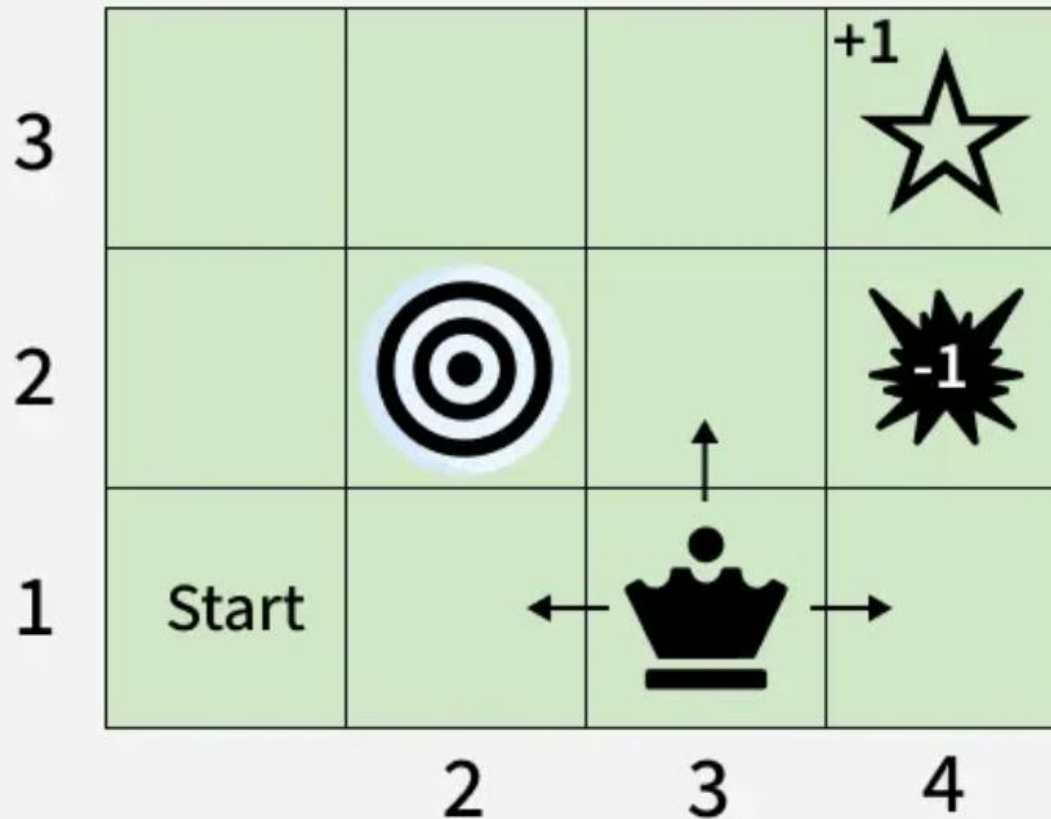
❑ **Reward (R):** A reward is a number given to the agent after it takes an action. If the reward is positive, it means the result of the action was good. If the reward is negative, it means the outcome was bad or there was a penalty help the agent learn what's good or bad.

Examples:

- +1 for reaching the goal
- -1 for stepping into fire
- -0.1 for each step to encourage fewer moves

❑ **Policy ( $\pi$ ):** A policy is the agent's plan. It tells the agent: “If you are in this state, take this action.” The goal is to find the best policy that helps the agent earn the highest total reward over time.

- ❑ Let's consider a 3x4 grid world. The agent starts at cell (1,1) and aims to reach the Blue Diamond at (4,3) while avoiding Fire at (4,2) and a Wall at (2,2). At each state the agent can take one of the following actions: UP, DOWN, LEFT or RIGHT



## ❑ **Movement with Uncertainty (Transition Model)**

- The agent's moves are stochastic (uncertain):
- 80% chance of going in the intended direction.
- 10% chance of going left of the intended direction.
- 10% chance of going right of the intended direction.



## ❑ Reward System

- +1 for reaching the goal.
- -1 for falling into fire.
- -0.04 for each regular move (to encourage shorter paths).
- 0 for hitting a wall (no movement or penalty).

## ❑ Goal and Policy

- The agent's objective is to maximize total rewards.
- It must find an optimal policy: the best action to take in each state to reach the goal quickly while avoiding danger.

## ❑ Path Example

- One possible optimal path is: UP  $\rightarrow$  UP  $\rightarrow$  RIGHT  $\rightarrow$  RIGHT  $\rightarrow$  RIGHT
- But because of randomness the agent must plan carefully to avoid accidentally slipping into fire.

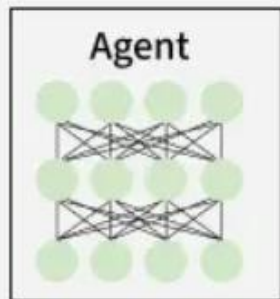
# Applications of MDP

- ❑ **Robots and Machines:** Robots use MDPs to decide how to move safely and efficiently in places like factories or warehouses and avoid obstacles.
- ❑ **Game Strategy:** In board games or video games MDPs help characters to choose the best moves to win or complete tasks even when outcomes are not certain.

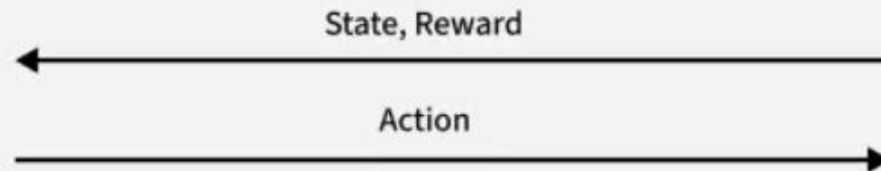
- ❑ **Healthcare:** Doctors can use it to plan treatments for patients, choosing actions that improve health while considering uncertain effects.
- ❑ **Traffic and Navigation:** Self-driving cars or delivery vehicles use it to find safe routes and avoid accidents on unpredictable roads.
- ❑ **Inventory Management:** Stores and warehouses use MDPs to decide when to order more stock so they don't run out or keep too much even when demand changes.

# Online vs. Offline RL

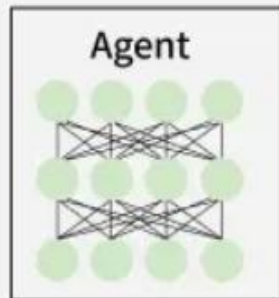
- ❑ Reinforcement Learning can be categorized based on how and when the learning agent acquires data from its environment, dividing the methods into online RL and offline RL (also known as batch RL).
- ❑ In Online RL, the agent learns by actively interacting with the environment in real-time. It collects fresh data during training by executing actions and observing immediate feedback as it learns.
- ❑ Offline RL trains the agent exclusively on a pre-collected static dataset of interactions generated by other agents, human demonstrations or historical logs. The agent does not interact with the environment during learning.



### Online reinforcement learning



**Environment**



### Offline reinforcement learning



**Environment**



<b>Aspect</b>	<b>Online RL</b>	<b>Offline RL</b>
<b>Data Acquisition</b>	Direct, real-time interaction with environment	Static, pre-collected dataset
<b>Adaptivity</b>	High, continuously adapts	Limited, depends on dataset coverage
<b>Suitability</b>	When environment access or simulation is feasible	When environment interaction is costly or risky
<b>Challenges</b>	Resource-intensive, potentially unsafe	Distributional shift, counterfactual inference issues



# Applications of Reinforcement Learning

- ❑ **Robotics:** RL is used to automate tasks in structured environments such as manufacturing, where robots learn to optimize movements and improve efficiency.
- ❑ **Games:** Advanced RL algorithms have been used to develop strategies for complex games like chess, Go and video games, outperforming human players in many instances.

- ❑ **Industrial Control:** RL helps in real-time adjustments and optimization of industrial operations, such as refining processes in the oil and gas industry.
- ❑ **Personalized Training Systems:** RL enables the customization of instructional content based on an individual's learning patterns, improving engagement and effectiveness.

# Advantages of Reinforcement Learning

- ❑ Solves complex sequential decision problems where other approaches fail.
- ❑ Learns from real-time interaction, enabling adaptation to changing environments.
- ❑ Does not require labeled data, unlike supervised learning.
- ❑ Can innovate by discovering new strategies beyond human intuition.
- ❑ Handles uncertainty and stochastic environments effectively.

# Disadvantages of Reinforcement Learning

- ❑ Computationally intensive, requiring large amounts of data and processing power.
- ❑ Reward function design is critical; poor design leads to unintended behaviors.
- ❑ Not suitable for simple problems where traditional methods are more efficient.
- ❑ Challenging to debug and interpret, making it hard to explain decisions.
- ❑ Exploration-exploitation trade-off requires careful balancing to optimize learning.