

Introduction to Reinforcement Learning

Naeemullah Khan
naeemullah.khan@kaust.edu.sa



جامعة الملك عبدالله
للعلوم والتكنولوجيا
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

June 30, 2025

Reinforcement Learning

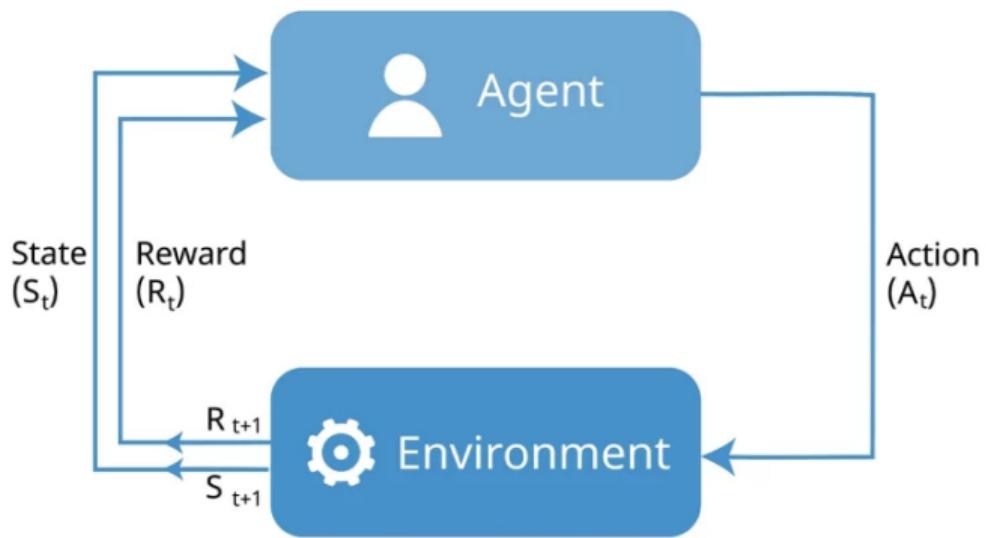


Table of Contents

1. Motivation
2. Learning Outcomes
3. Introduction
4. Types of Reinforcement Learning
5. Markov Decision Processes (MDPs)
6. Policy
7. Value Function
8. Q-Value Function
9. Bellman Equation
 1. Optimal Bellman Equation
10. Q-Learning
11. Limitations
12. Summary

Table of Contents (cont.)

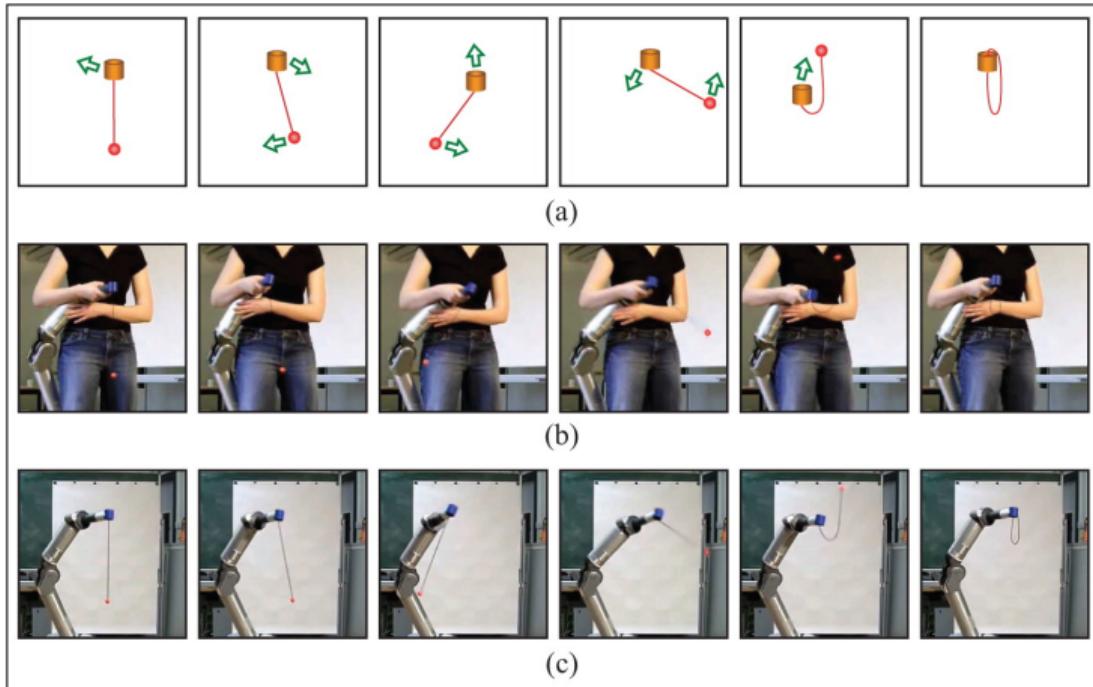
13. References

Motivation



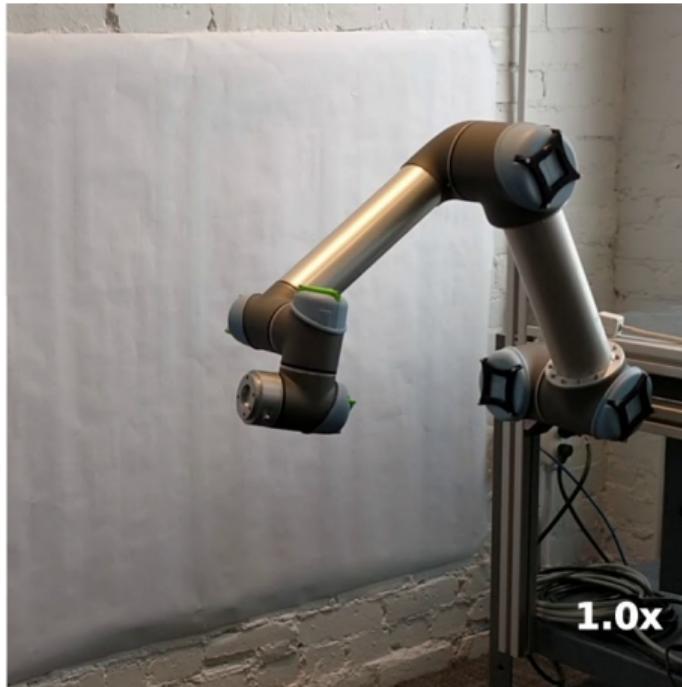
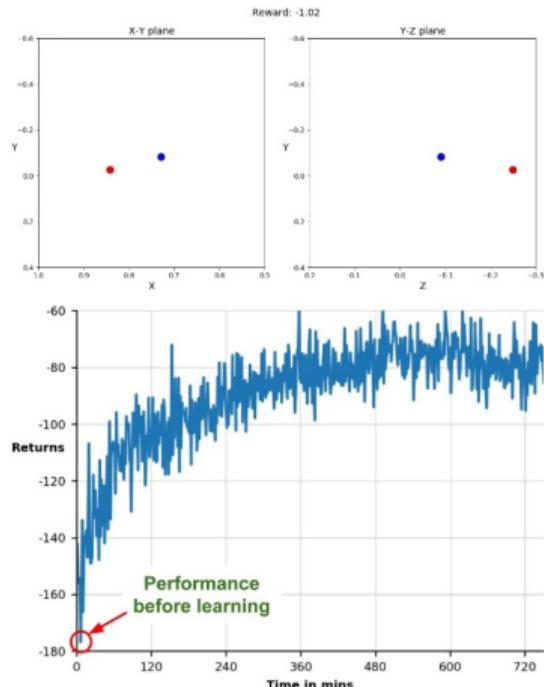
Cart-Pole Swing-Up Task

Motivation (cont.)



Robot Ball-in-a-Cup Task (Learning Motor Primitives in Robotics)

Motivation (cont.)

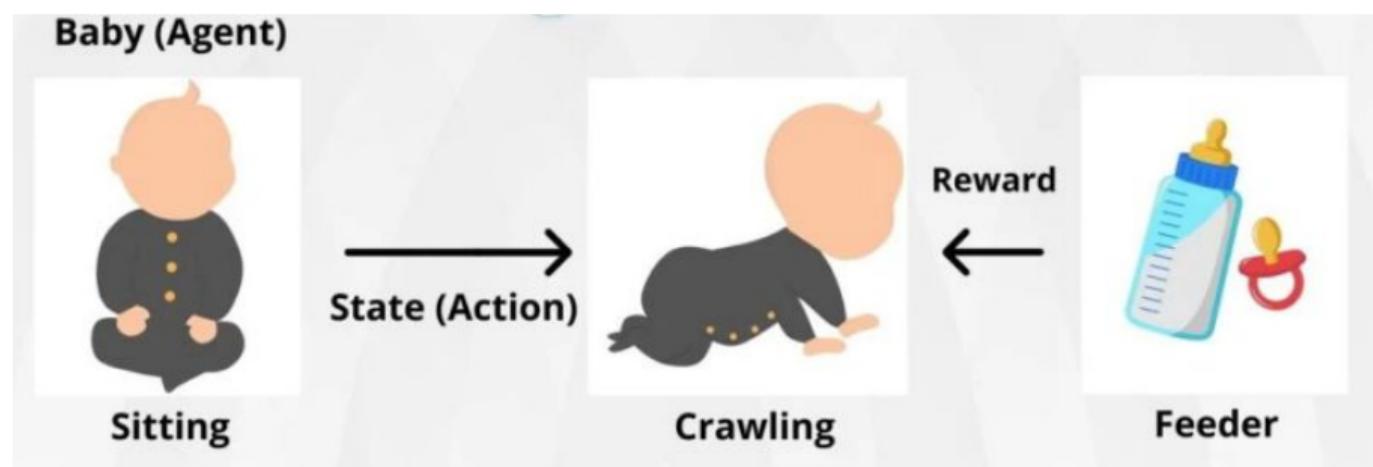


Robot Learning in Real-World Environments

Motivation (cont.)

- ▶ Many AI problems involve **sequential decision-making** under uncertainty.
- ▶ Inspired by **animal learning** (reward-based adaptation).
- ▶ Reinforcement Learning (RL) allows agents to learn optimal behaviors by **trial and error**.
- ▶ Real-world success: AlphaGo, Robotics, Game AI, Recommendation Systems.

Motivation (cont.)

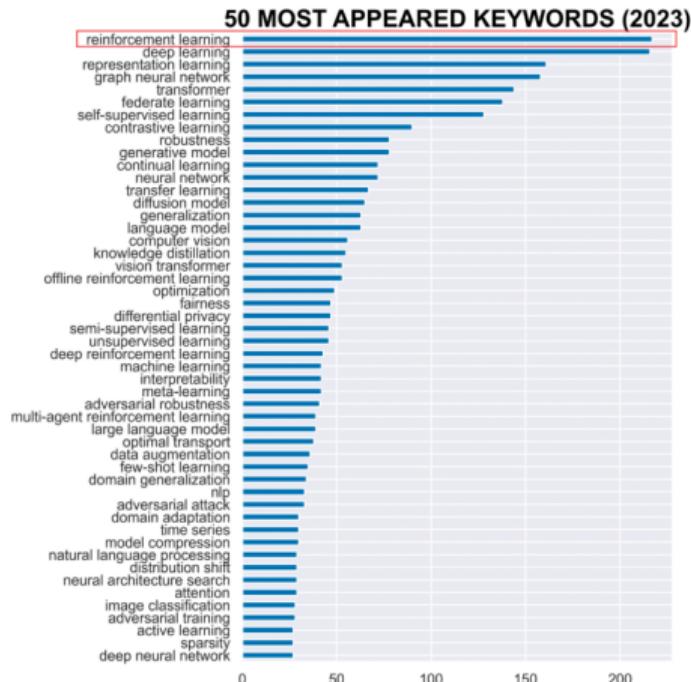


Baby Learning to Crawl towards a Milk Bottle (Reward)

By the end of this session, you will be able to:

- ▶ Understand the fundamental **principles and structure of RL**.
- ▶ Model problems using **Markov Decision Processes (MDPs)**.
- ▶ Implement and apply **Q-learning** to solve control problems.
- ▶ Identify **practical applications and challenges** in RL.
- ▶ Critically analyze the **limitations and future trends** of RL.

Reinforcement Learning: **Introduction**



Source: ICLR 2023 - Open Review Data

Introduction

Definition: RL is a type of machine learning where agents learn to make decisions by interacting with an environment.

Core Elements:

- ▶ **Agent**: The learner or decision maker.
- ▶ **Environment**: The external system the agent interacts with.
- ▶ **State**: A representation of the current situation.
- ▶ **Action**: Choices the agent can make.
- ▶ **Reward**: Feedback signal for evaluating actions.

Introduction (cont.)



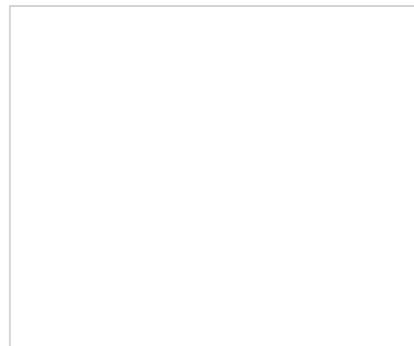
agent



Introduction (cont.)



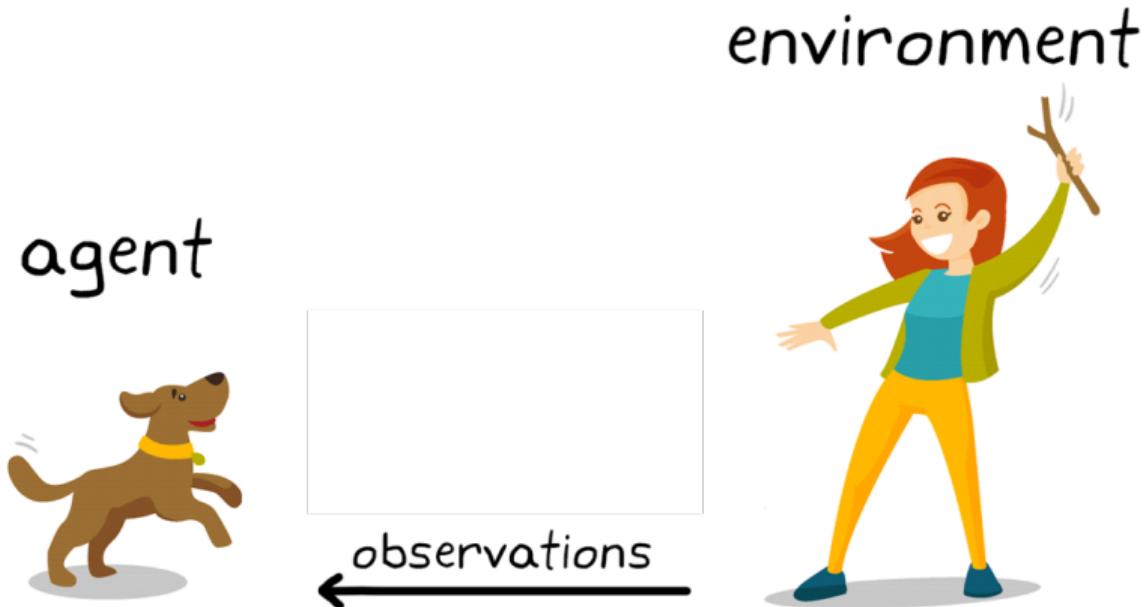
agent



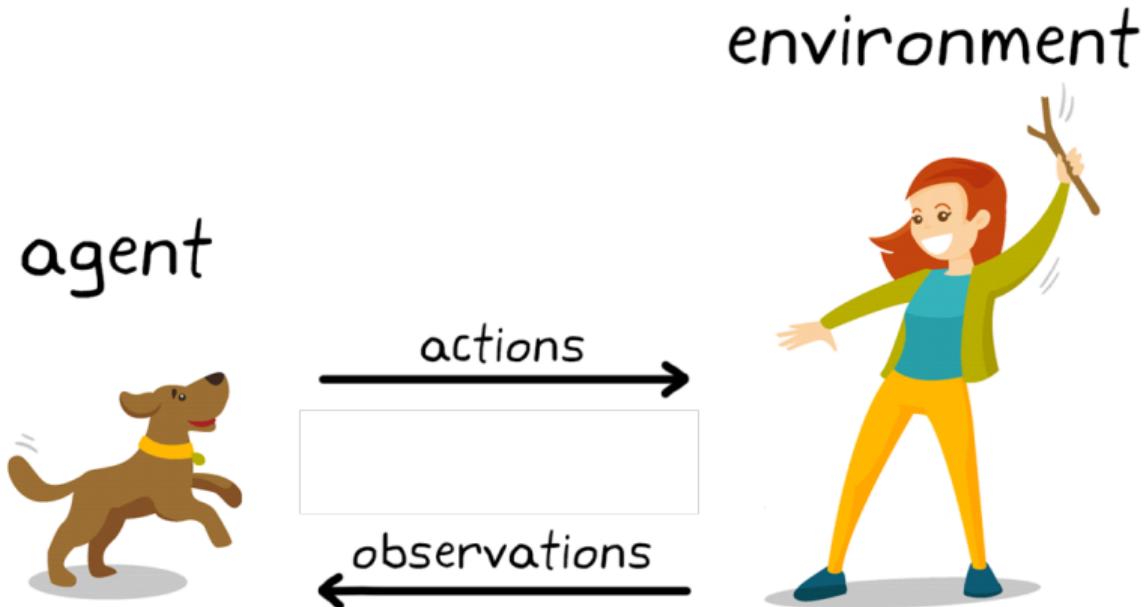
environment



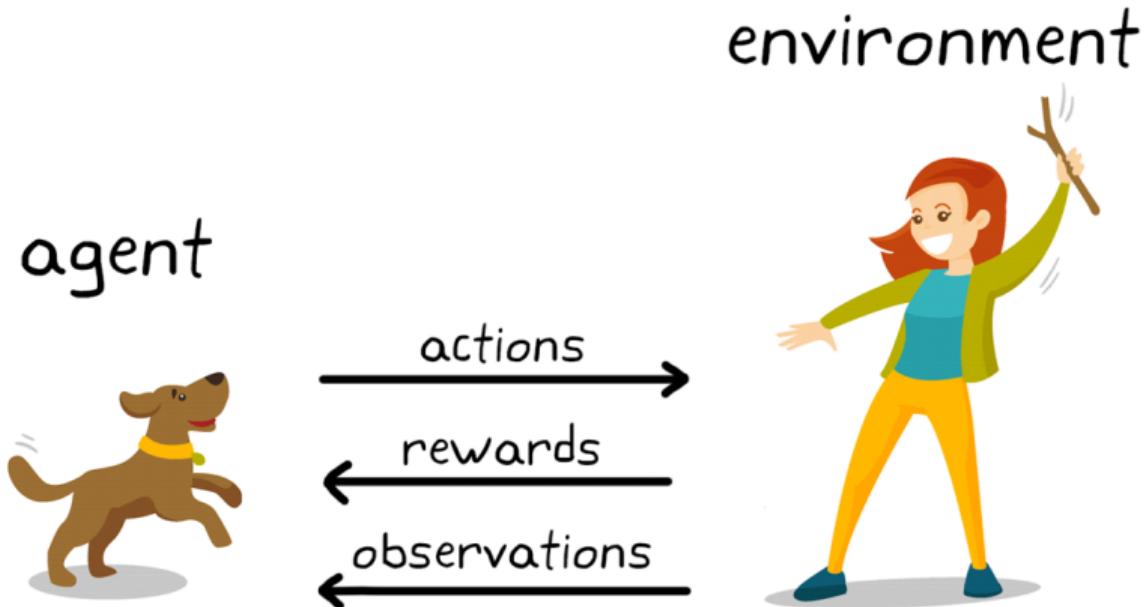
Introduction (cont.)



Introduction (cont.)



Introduction (cont.)

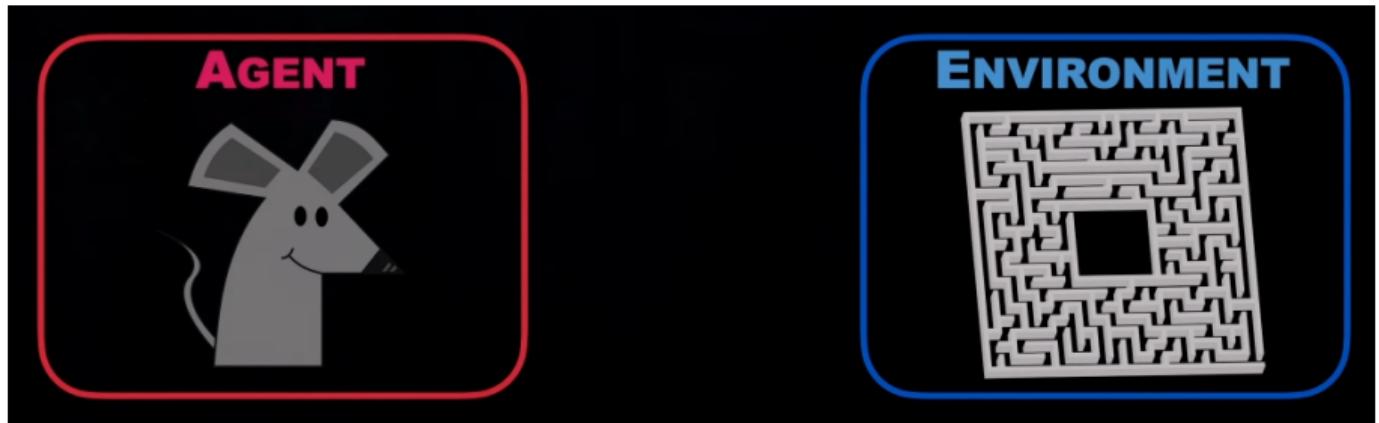


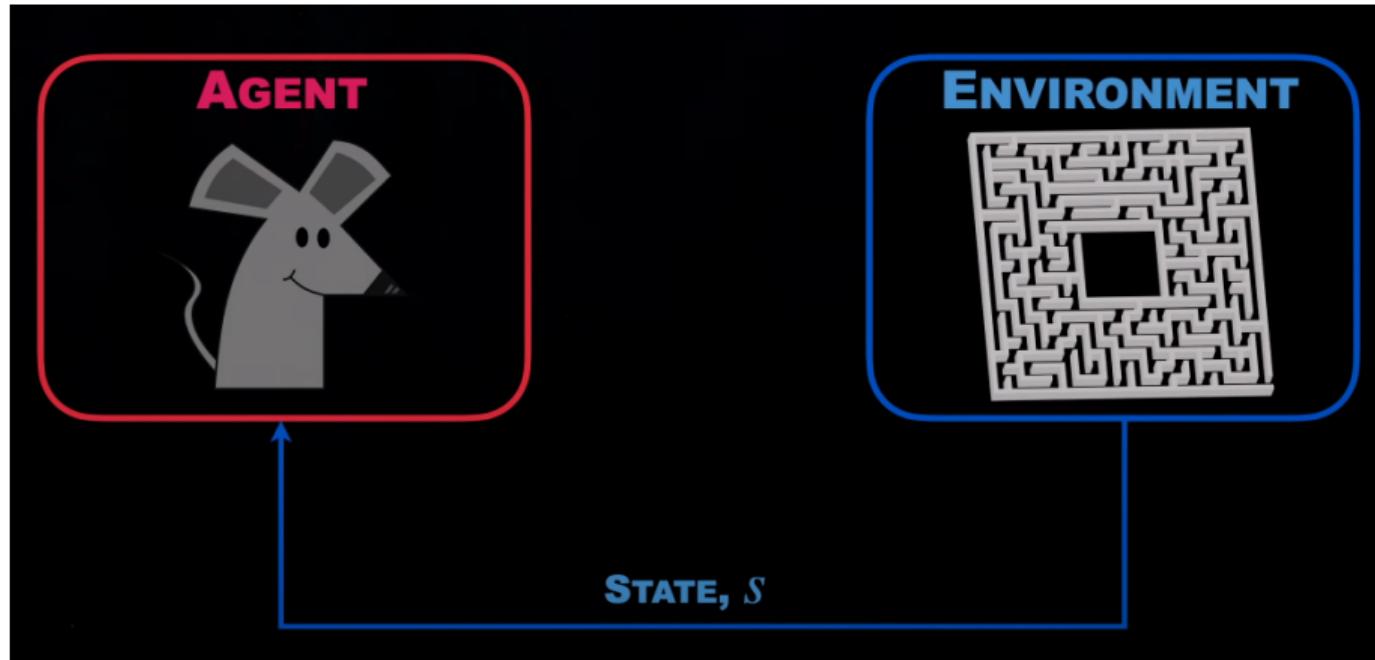
Key Concepts:

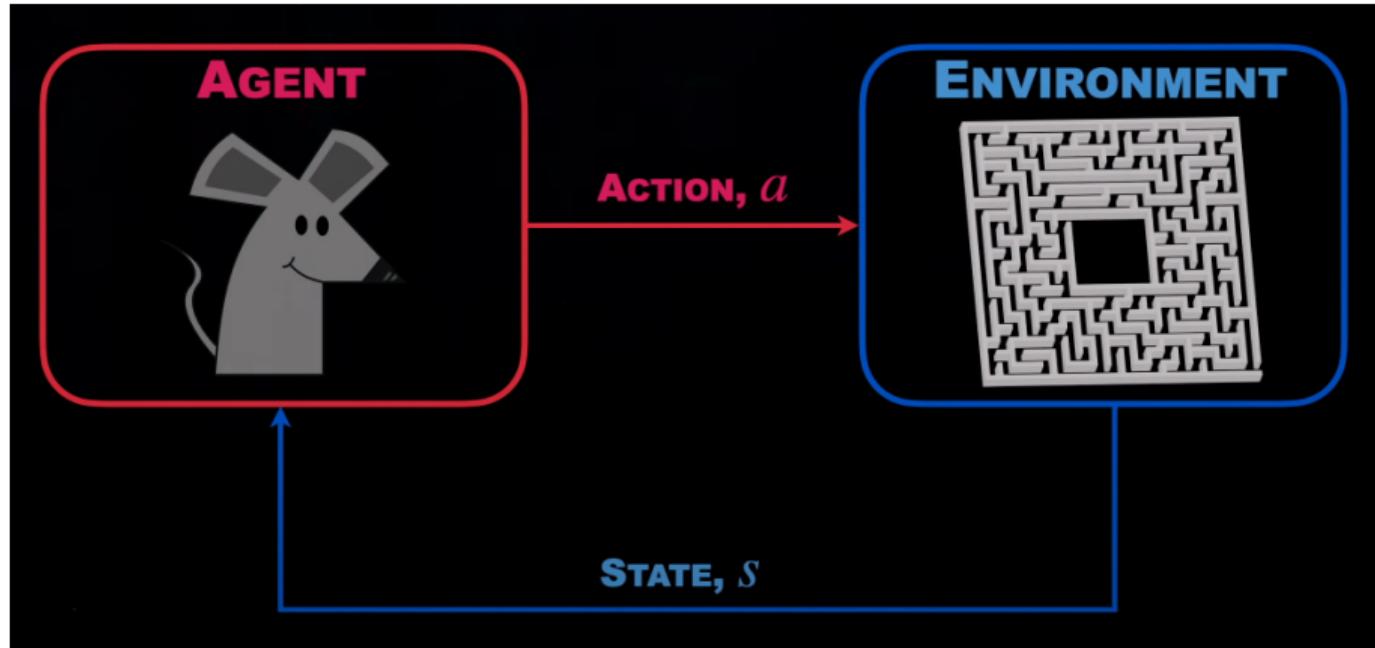
- ▶ Agent interacts with environment.
- ▶ Learns by trial and error.

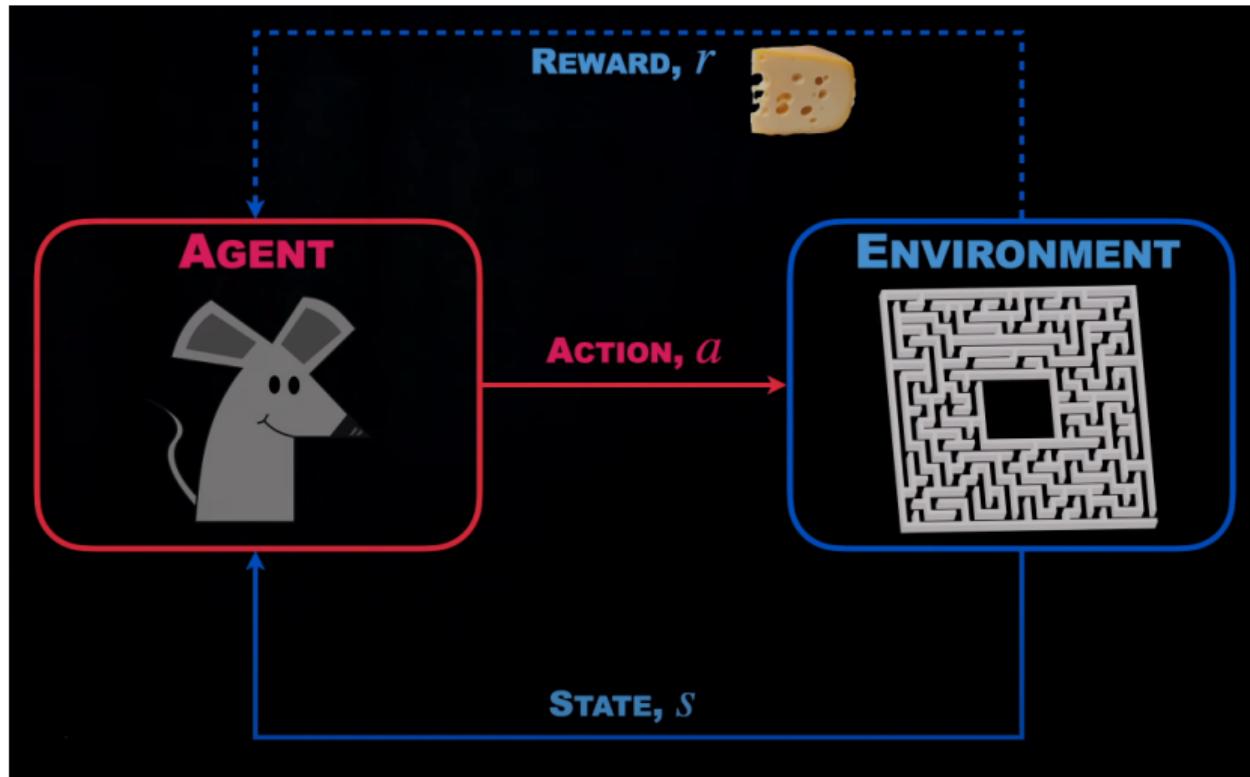
How does it work?

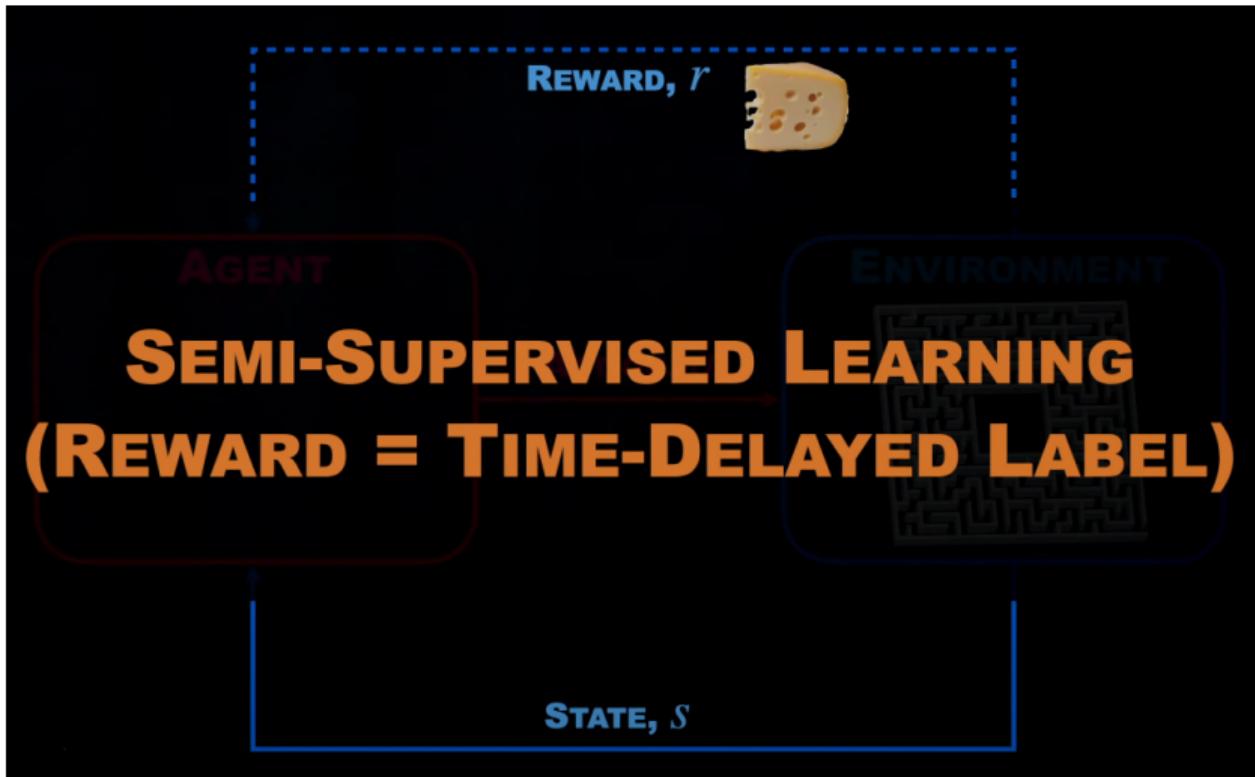
- ▶ Agent gets rewards or penalties.
- ▶ Goal: Maximize cumulative reward (a.k.a. return).

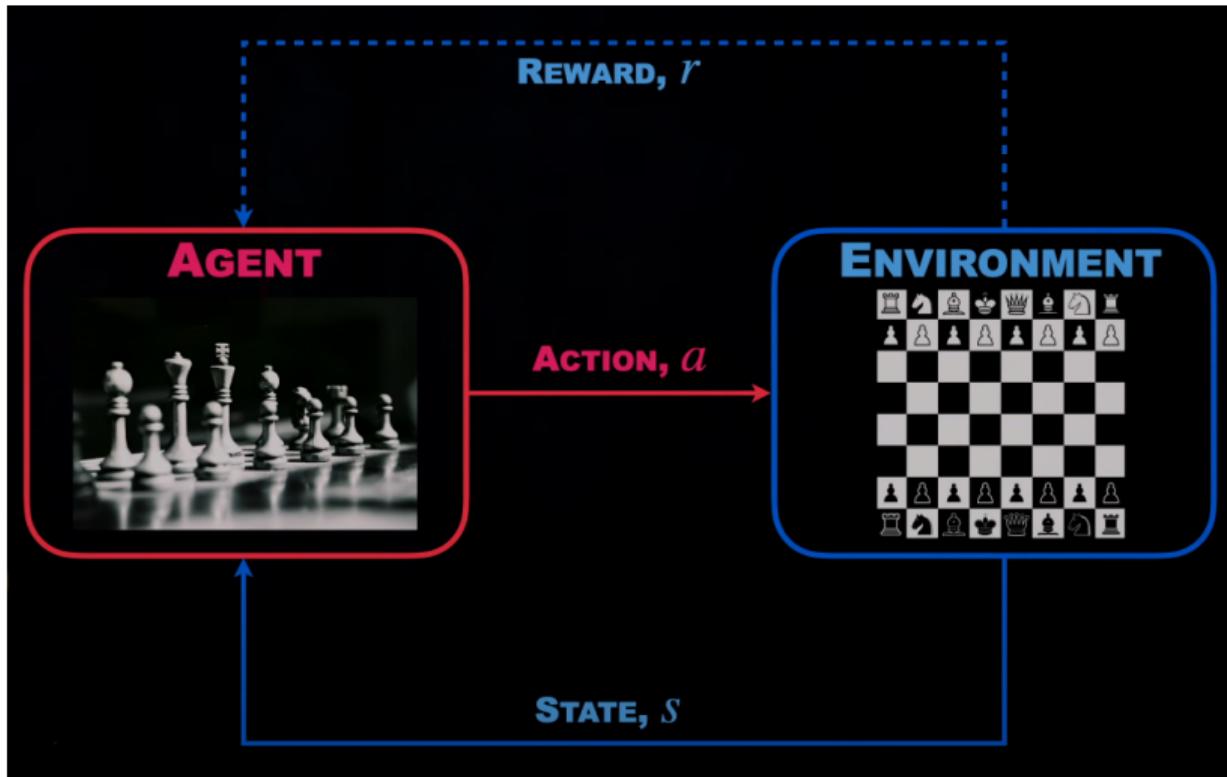


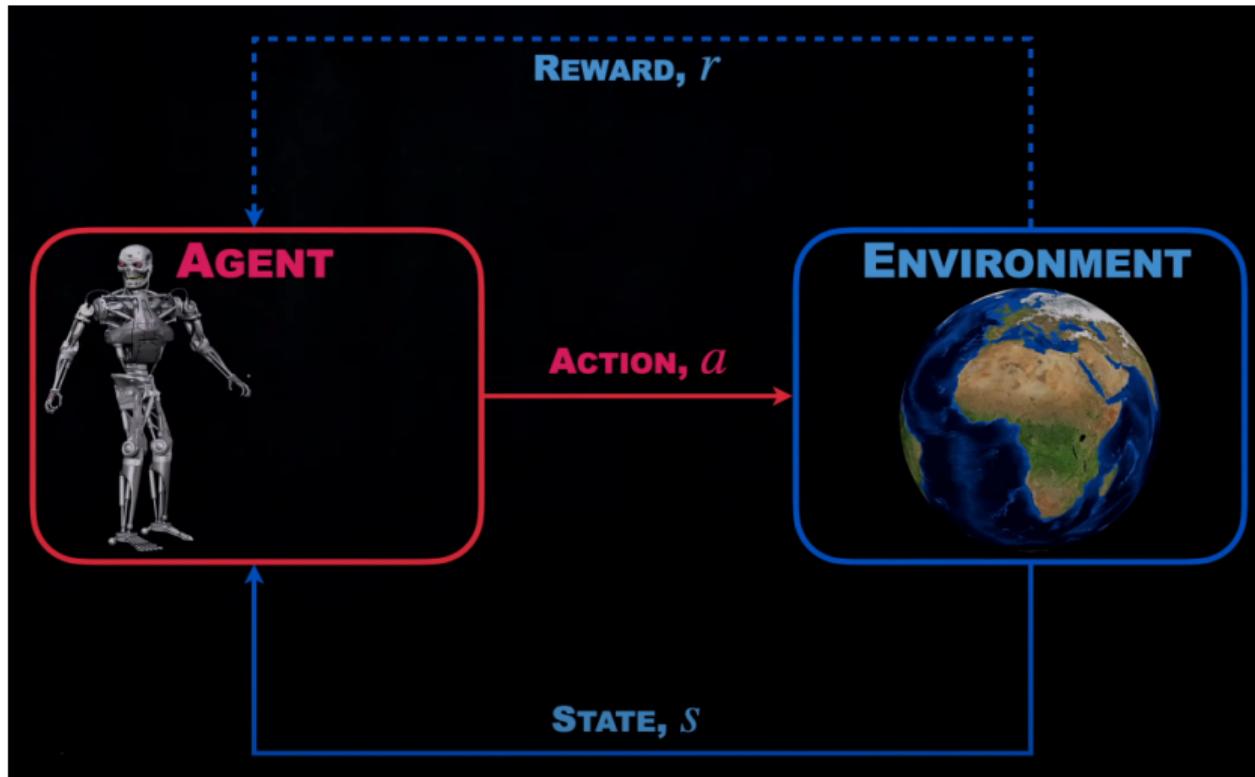


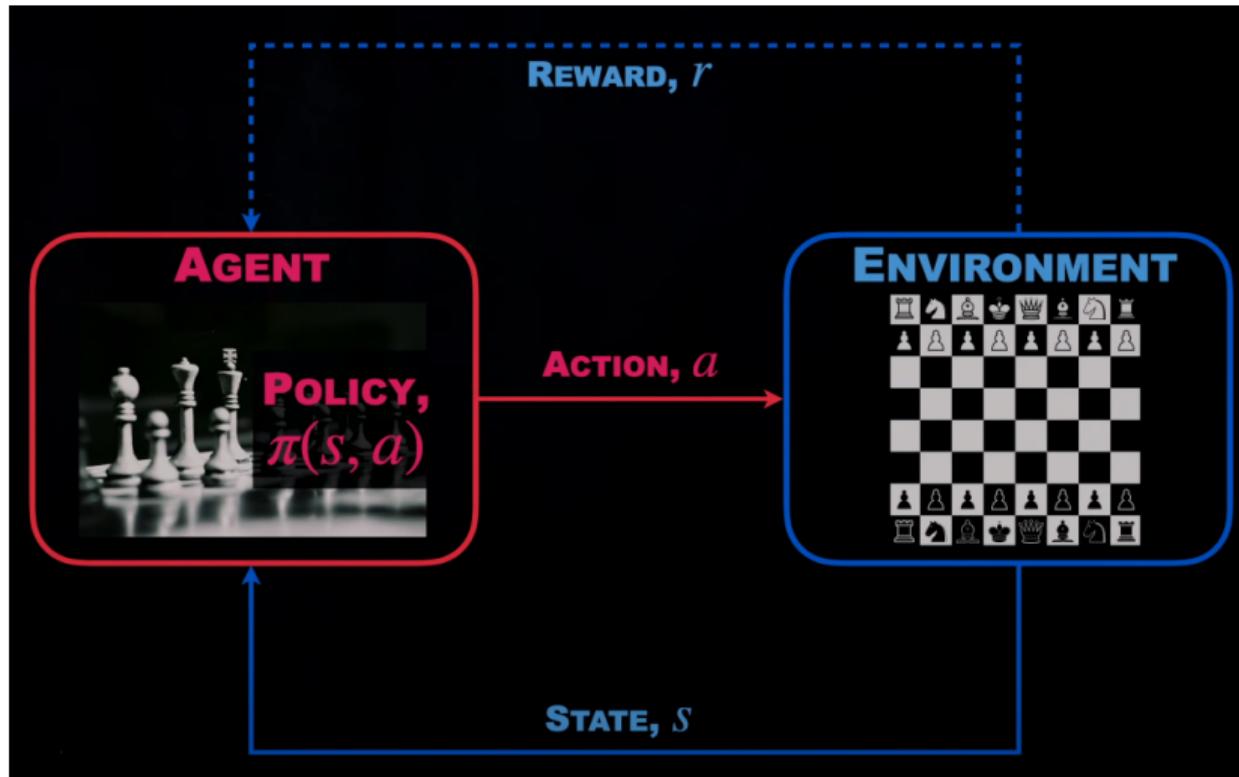












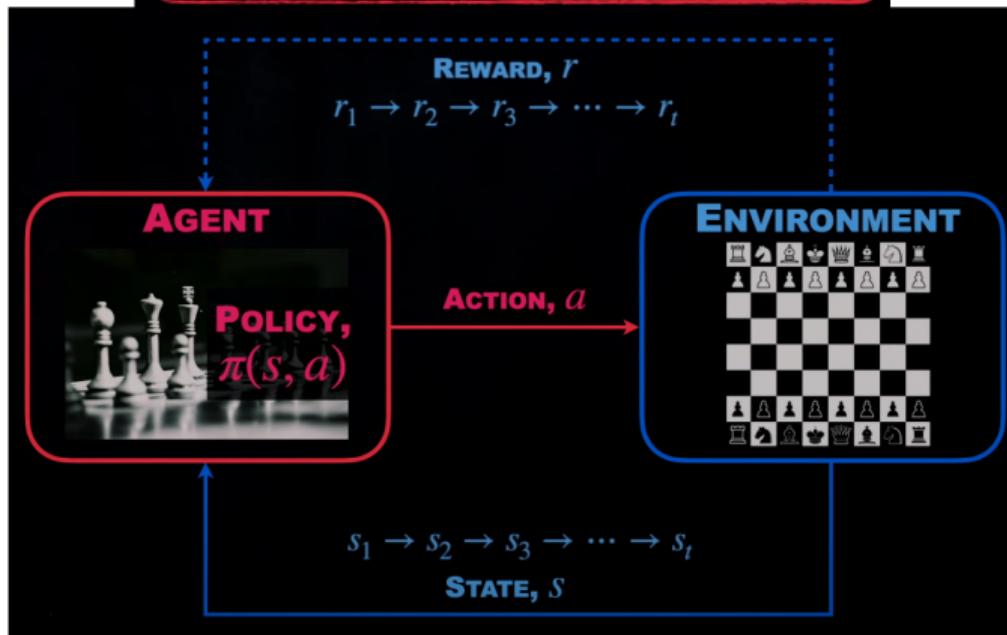


POLICY $\pi(s, a) = \Pr(a = a | s = s)$



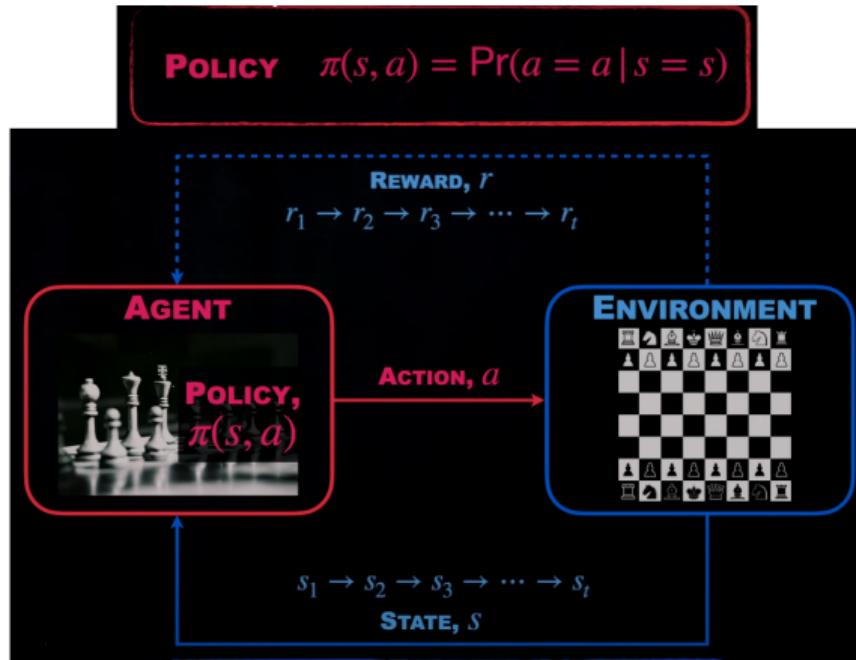


POLICY $\pi(s, a) = \Pr(a = a | s = s)$





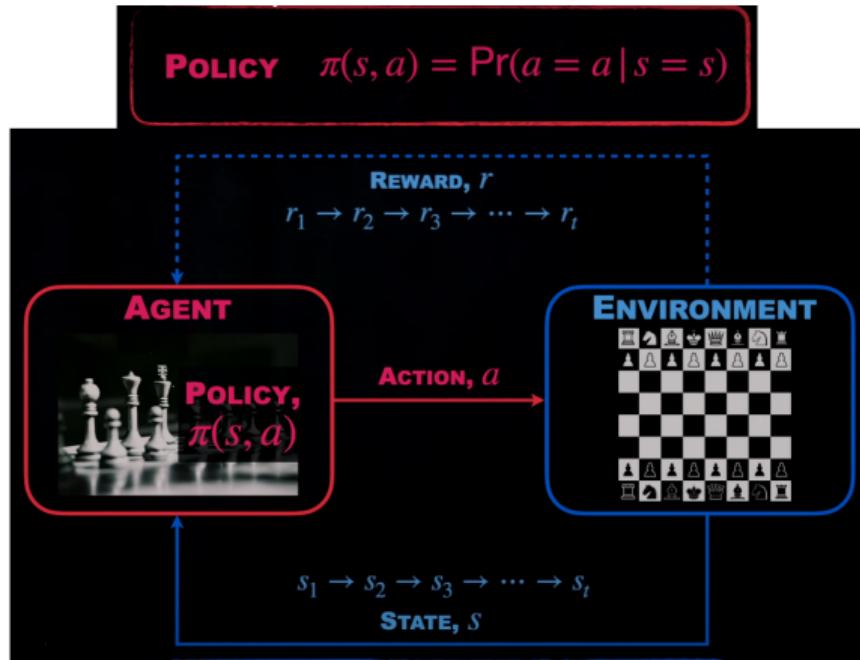
$$\textbf{POLICY} \quad \pi(s, a) = \Pr(a = a | s = s)$$



$$\textbf{VALUE} \quad V_{\pi}(s) = \mathbb{E} \left(\sum_t \gamma^t r_t | s_0 = s \right)$$



$$\text{POLICY} \quad \pi(s, a) = \Pr(a = a | s = s)$$



$$\text{VALUE} \quad V_{\pi}(s) = \mathbb{E} \left(\sum_t \gamma^t r_t | s_0 = s \right)$$

Discount Rate



POLICY $\pi(s, a) = \Pr(a = a | s = s)$

REWARD, r

$r_1 \rightarrow r_2 \rightarrow r_3 \rightarrow \dots \rightarrow r_t$

AGENT

**GOAL: OPTIMIZE POLICY TO
MAXIMIZE FUTURE REWARDS**

ENVIRONMENT

$s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow \dots \rightarrow s_t$

STATE, s

VALUE $V_\pi(s) = \mathbb{E} \left(\sum_{t=0}^T r_t | s_0 = s \right)$

DISCOUNT RATE

Supervised Learning vs. Reinforcement Learning

Supervised Learning

- ▶ Given labeled data: $\{(x_i, y_i)\}$, learn $f(x) \approx y$
- ▶ Directly told what to output
- ▶ Inputs x are independently, identically distributed (i.i.d.)

Reinforcement Learning

- ▶ Learn behavior $\pi(a|s)$
- ▶ Learn from experience, indirect feedback
- ▶ Data not i.i.d.: actions affect future observations

Examples of Reinforcement Learning

ChatGPT

Step 1

Collect demonstration data and train a supervised policy.

A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3.5 with supervised learning.



Step 2

Collect comparison data and train a reward model.

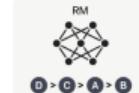
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.



Write a story about otters.



The PPO model is initialized from the supervised policy.

The policy generates an output.



The reward model calculates a reward for the output.



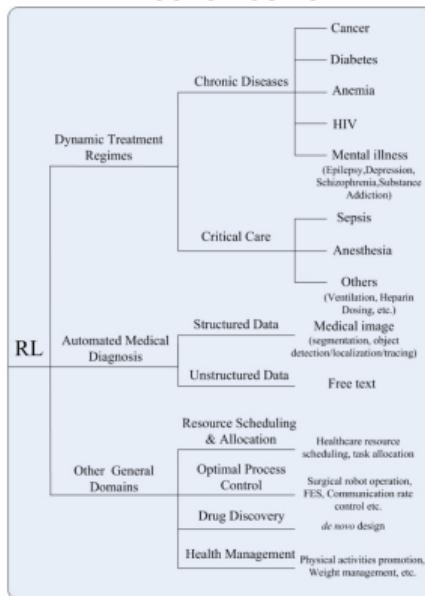
The reward is used to update the policy using PPO.

r_k

<https://openai.com/blog/chatgpt>

Examples of Reinforcement Learning (cont.)

Healthcare

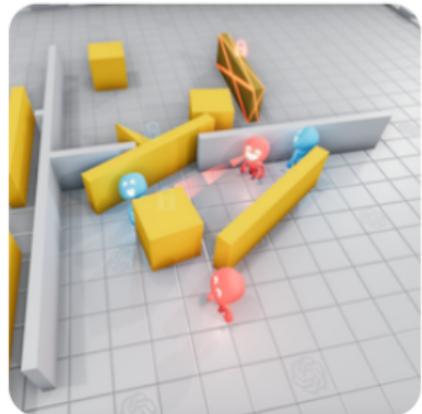


Yu, et al. Reinforcement Learning in Healthcare: A Survey

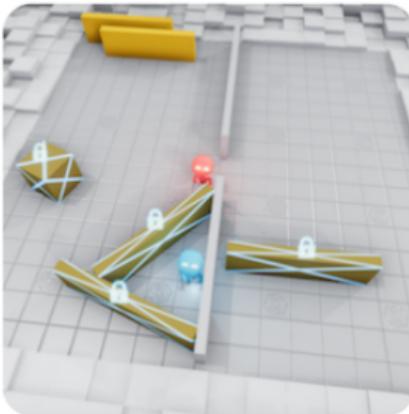
Examples of Reinforcement Learning (cont.)

Multiagent Hide & Seek

(a) Running and Chasing



(b) Fort Building

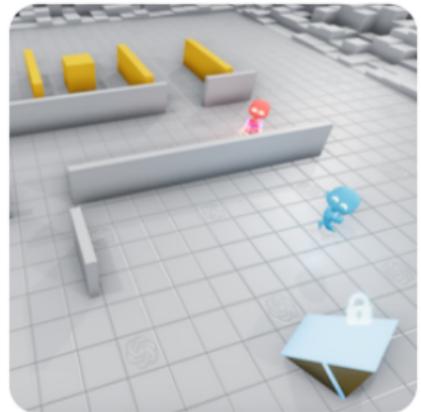


(c) Ramp Use



Examples of Reinforcement Learning (cont.)

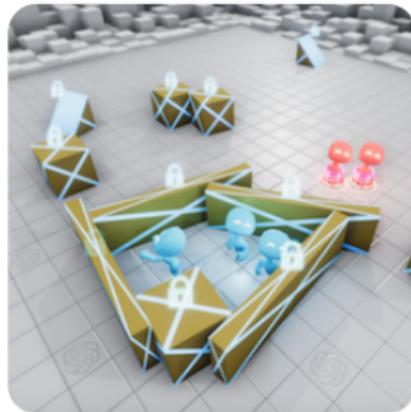
Multiagent Hide & Seek



(d) Ramp Defense



(e) Box Surfing



(f) Surf Defense

Examples of Reinforcement Learning (cont.)

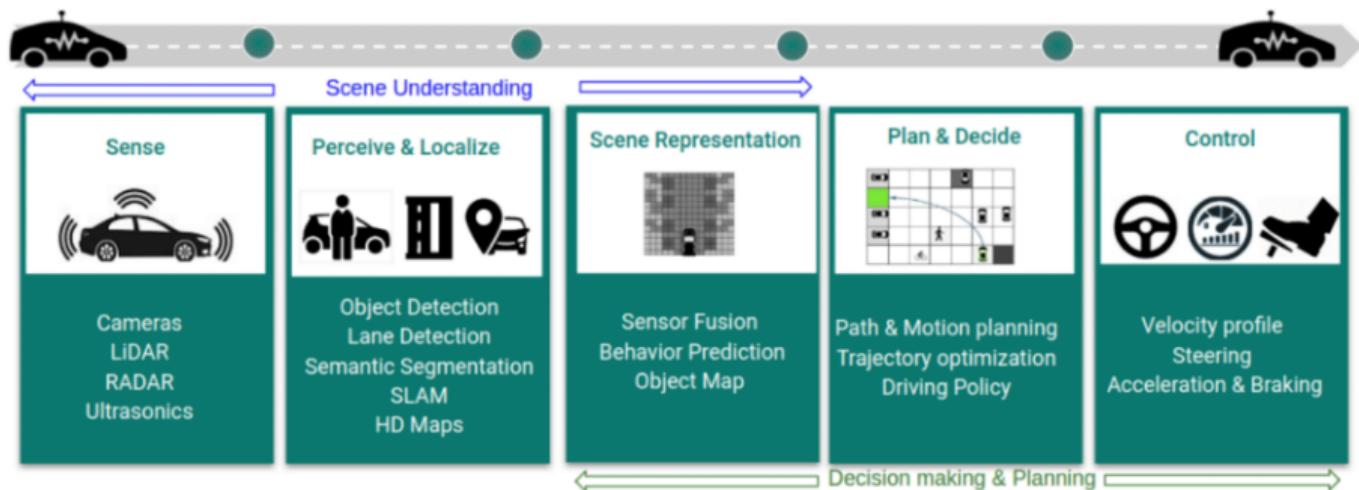
Gran Turismo GTSophie



<https://www.gran-turismo.com/us/gran-turismo-sophy/technology/>

Examples of Reinforcement Learning (cont.)

Autonomous Driving



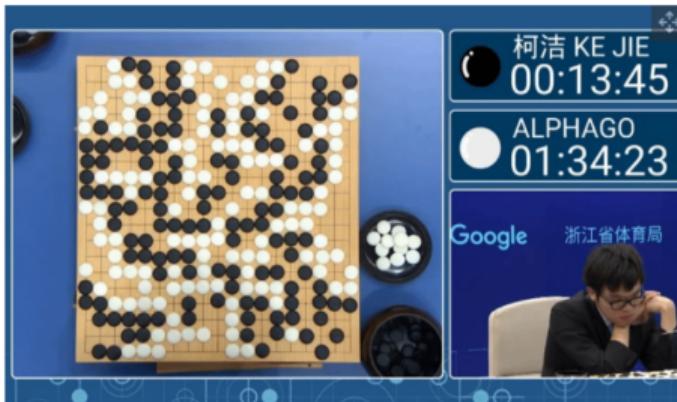
Kiran, et al. Deep Reinforcement Learning for Autonomous Driving: A Survey

Examples of Reinforcement Learning (cont.)

AlphaGo

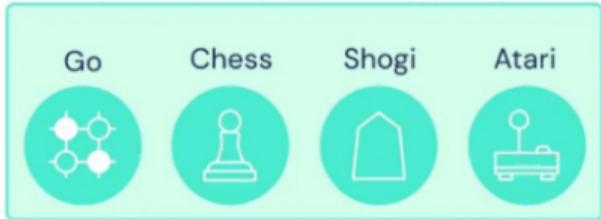
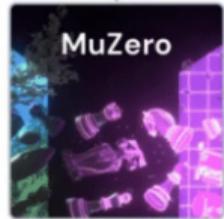
'AlphaGo' AI Scores Narrow Win Against Ke Jie, World's Top 'Go' Player

By Lucian Armasu published May 23, 2017



<https://www.tomshardware.com/news/alphago-narrow-win-ke-jie,34486.html>

Examples of Reinforcement Learning (cont.)



MuZero learns the rules of the game, allowing it to also master environments with unknown dynamics.
(Dec 2020, Nature)

<https://www.deepmind.com/blog/muzero-mastering-go-chess-shogi-and-atari-without-rules>

Examples of Reinforcement Learning (cont.)

Minecraft

AI learns how to play Minecraft by watching videos



Open AI has trained a neural network to play Minecraft by Video PreTraining (VPT) on a massive unlabeled video dataset of human Minecraft play, while using just a small amount of labeled contractor data.

With a bit of fine-tuning, the AI research and deployment company is confident that its model can learn to craft diamond tools, a task that usually takes proficient humans over 20 minutes (24,000 actions). Its model uses the native human interface of keypresses and mouse...



29 June 2022 | Artificial Intelligence

<https://www.artificialintelligence-news.com/2022/06/29/ai-learns-how-to-play-minecraft-by-watching-videos/>

Examples of Reinforcement Learning (cont.)

More Examples



BIOTECHNOLOGY

Reinforcement learning: From board games to protein design

Scientists have successfully applied reinforcement learning to a challenge in molecular biology. The team of researchers developed powerful new protein design software adapted from a strategy proven adept at board games like ...



APR 20, 2023

0

37



OPTICS & PHOTONICS

Chiral detection of biomolecules based on reinforcement learning

As one of the basic physical properties, chirality plays an important role in many fields. Especially in biomedical chemistry, the discrimination of enantiomers is a very important research subject. Most biomolecules exhibit

...



FEB 22, 2023

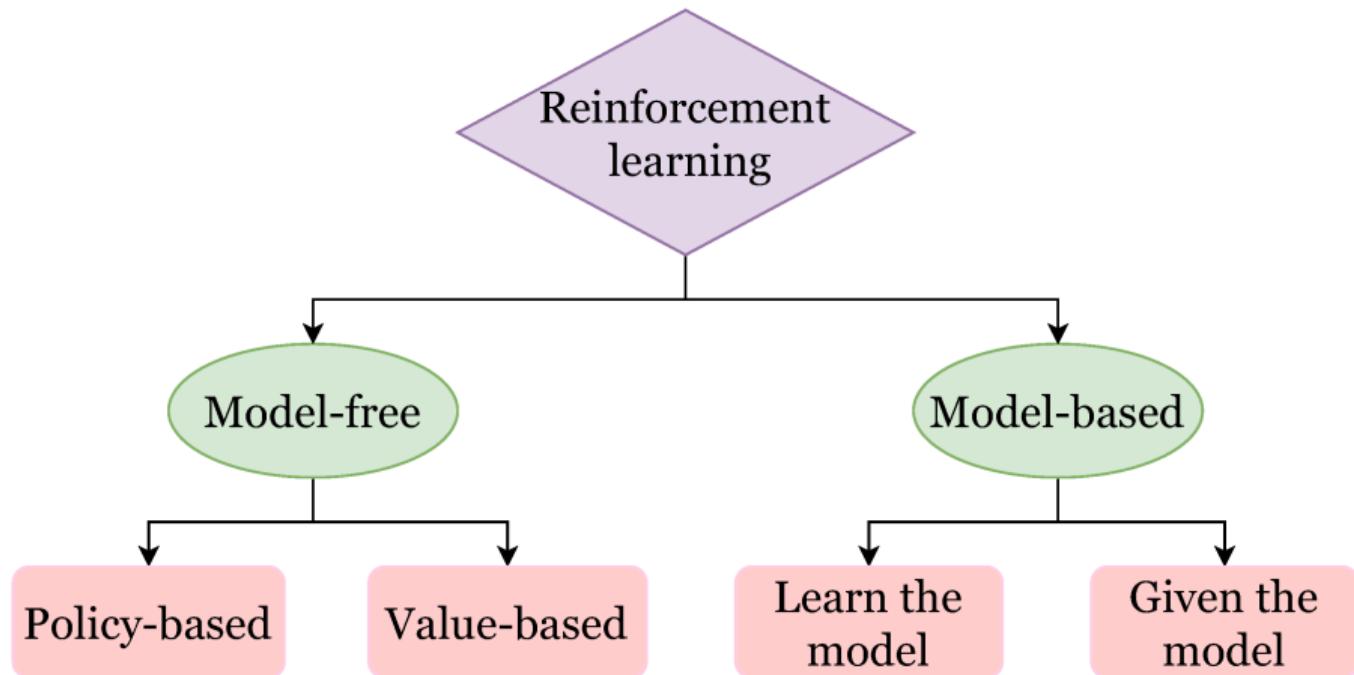
0

4

<https://phys.org/tags/reinforcement+learning/>

Reinforcement Learning: Types

Types of Reinforcement Learning



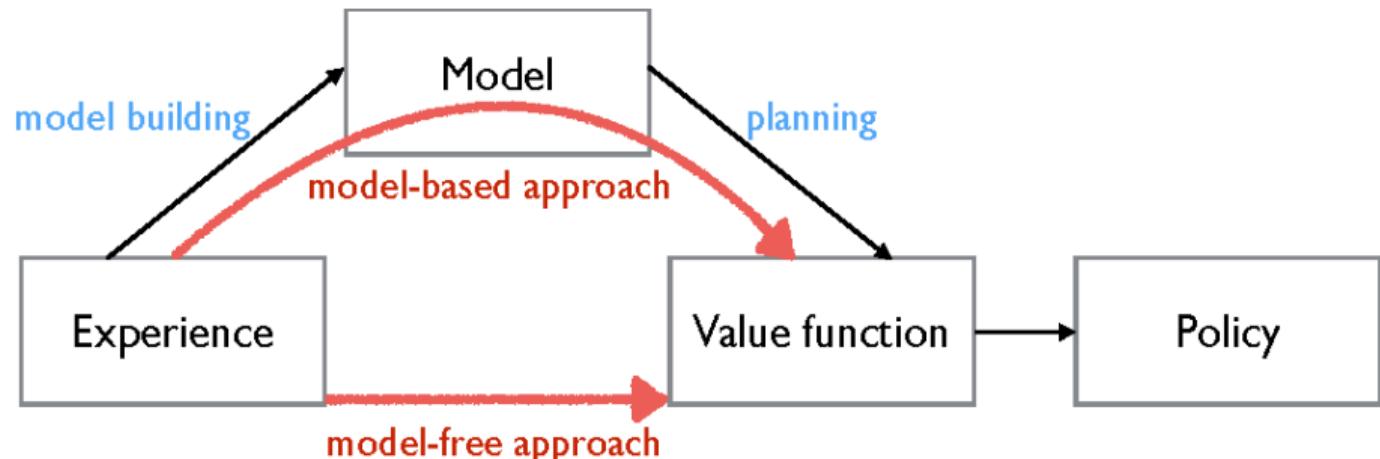
Model-based vs. Model-free

- ▶ **Model-based:** The agent builds or uses a model of the environment to plan actions.
- ▶ **Model-free:** The agent learns to act without explicitly modeling the environment.

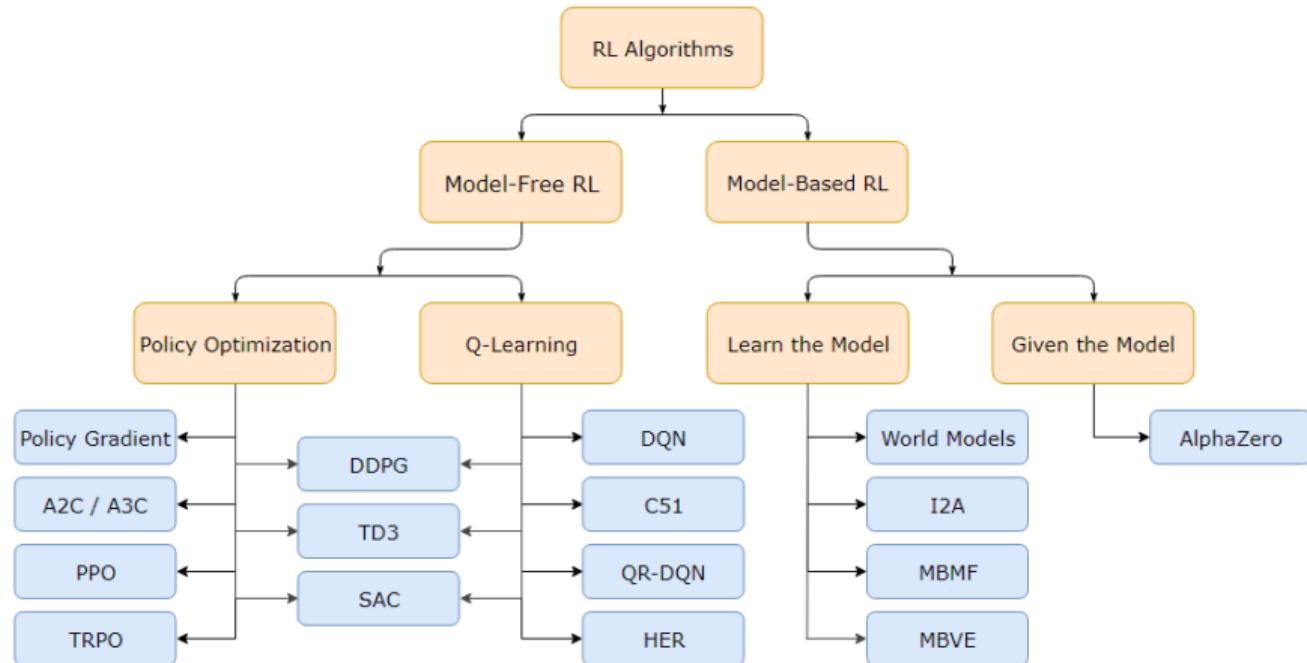
Value-based vs. Policy-based

- ▶ **Value-based:** The agent learns a value function (e.g., Q-learning) to evaluate actions or states.
- ▶ **Policy-based:** The agent directly learns a policy that maps states to actions (e.g., Policy Gradient methods).

Types of Reinforcement Learning (cont.)



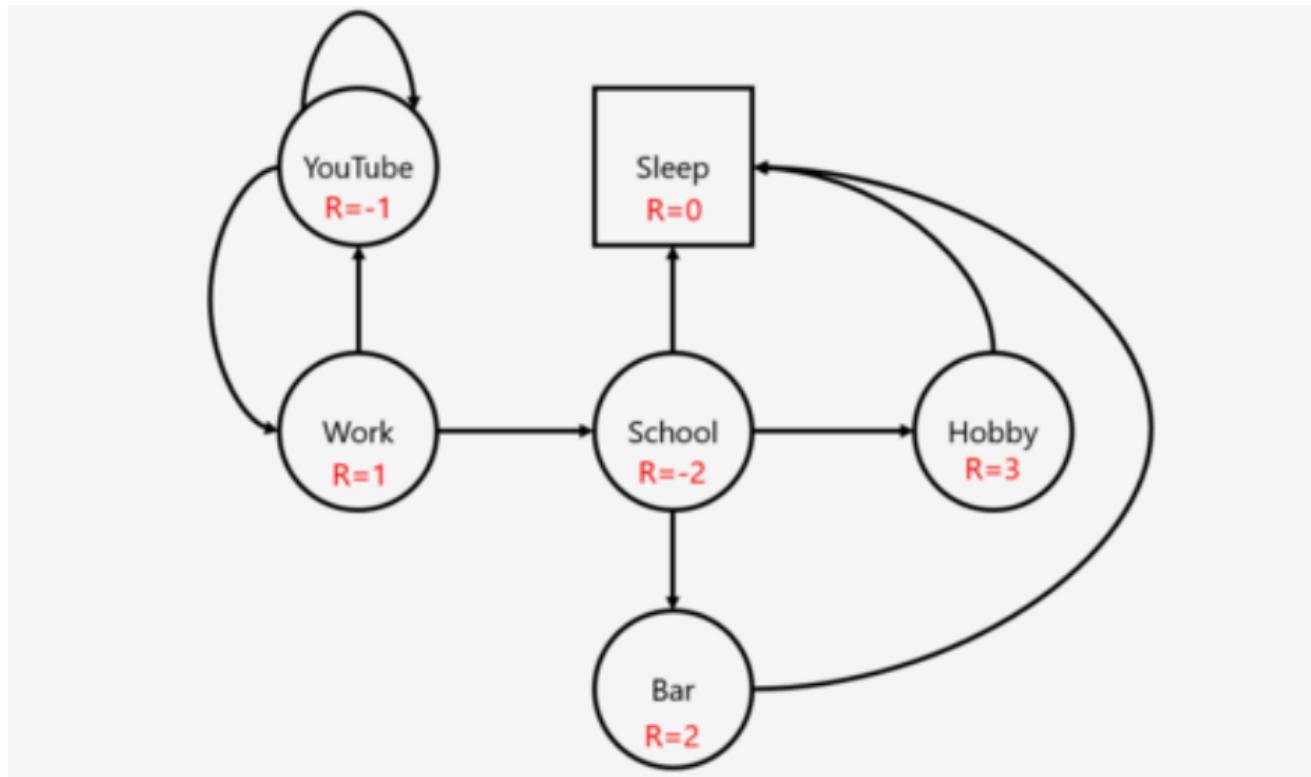
Types of Reinforcement Learning (cont.)



Reinforcement Learning: Markov Decision Processes (MDPs)

- ▶ **Markov property:** The current state completely characterizes the state of the world; that is, the future is independent of the past given the present.
- ▶ **Defined by:** $(\mathcal{S}, \mathcal{A}, R, P, \gamma)$
 - \mathcal{S} : set of possible states
 - \mathcal{A} : set of possible actions
 - $R(s, a)$: reward distribution given a (state, action) pair
 - $P(s' | s, a)$: transition probability, i.e., distribution over next state given a (state, action) pair
 - γ : discount factor
- ▶ Assumes the **Markov property**: the next state depends only on the current state and action.

Mathematical Formulation of the RL Problem (cont.)



Mathematical Formulation of the RL Problem (cont.)

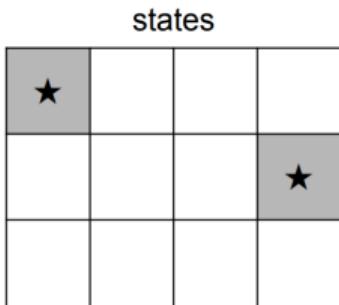
- ▶ At time step $t = 0$, the environment samples the initial state $s_0 \sim p(s_0)$.
- ▶ Then, for $t = 0$ until termination:
 - The agent selects action a_t .
 - The environment samples reward $r_t \sim R(\cdot|s_t, a_t)$.
 - The environment samples the next state $s_{t+1} \sim P(\cdot|s_t, a_t)$.
 - The agent receives reward r_t and next state s_{t+1} .
- ▶ A policy π is a function from \mathcal{S} to \mathcal{A} that specifies which action to take in each state.
- ▶ **Objective:** Find a policy π^* that maximizes the cumulative discounted reward: $\sum_{t \geq 0} \gamma^t r_t$.

Mathematical Formulation of the RL Problem (cont.)

- ▶ The agent and the environment interact in a time-sequenced loop:
 - The agent observes the current state and reward, then selects an action.
 - The environment responds by producing the next state and a scalar reward.
- ▶ Each state satisfies the **Markov property**:
 - The next state and reward depend only on the current state and action.
 - The current state captures all relevant information from the history.
 - The current state is a sufficient statistic for predicting the future (given the action).
- ▶ The goal of the agent is to maximize the expected sum of all future rewards by controlling the policy function $\pi : \mathcal{S} \rightarrow \mathcal{A}$.
- ▶ This setup forms a dynamic, time-sequenced control system under uncertainty.

A Simple MDP: Grid World

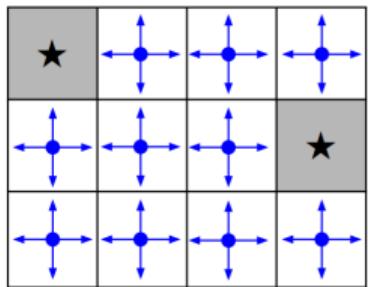
```
actions = {  
    1. right ←→  
    2. left ←→  
    3. up ↑↓  
    4. down ↑↓  
}
```



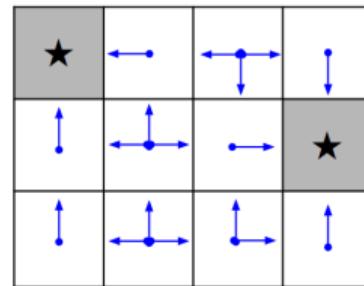
Set a negative “reward”
for each transition
(e.g. $r = -1$)

Objective: Reach one of the terminal states (greyed out) in the least number of actions.

A Simple MDP: Grid World (cont.)



Random Policy



Optimal Policy

Real-World Problems that Fit the MDP Framework

- ▶ Self-driving vehicles: choosing speed and steering to optimize safety and travel time
- ▶ Playing chess: receiving a reward only at the end of the game
- ▶ Complex logistical operations, such as managing movements in a warehouse
- ▶ Enabling a humanoid robot to walk or run on challenging terrains
- ▶ Managing an investment portfolio over time
- ▶ Controlling a power station efficiently
- ▶ Making optimal decisions during a football match
- ▶ Developing strategies to win an election (a high-complexity MDP)

Why Are These Problems Hard?

- ▶ The state space can be large or complex (involving many variables).
- ▶ Sometimes, the action space is also large or complex.
- ▶ There is no direct feedback on “correct” actions (the only feedback is the reward).
- ▶ Time-sequenced complexity: actions influence future states and actions.
- ▶ Actions can have delayed consequences (delayed rewards).
- ▶ The agent often does not know the model of the environment.
- ▶ The “model” refers to the probabilities of state transitions and rewards.
- ▶ Thus, the agent has to learn the model and solve for the optimal policy.
- ▶ Agent actions need to trade off between “exploration” and “exploitation”.

Reinforcement Learning: Policy

- ▶ Policy π determines how the agent chooses actions
- ▶ $\pi : S \rightarrow A$, mapping from states to actions
- ▶ Deterministic Policy:

$$\pi(s) = a$$

- ▶ Stochastic Policy:

$$\pi(a|s) = Pr(a_t = a | s_t = s)$$

The Optimal Policy

- ▶ We want to find optimal policy π^* that maximizes the sum of reward
- ▶ But, how do we handle the randomness (initial state, transition probability...)?

The Optimal Policy

- ▶ We want to find optimal policy π^* that maximizes the sum of reward
- ▶ But, how do we handle the randomness (initial state, transition probability...)?
- ▶ **Solution:** Maximize the expected sum of rewards!
- ▶ Formally,

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi \right] \quad \text{with} \quad s_0 \sim p(s_0), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$$

Reinforcement Learning: **Value Function**

Value Function

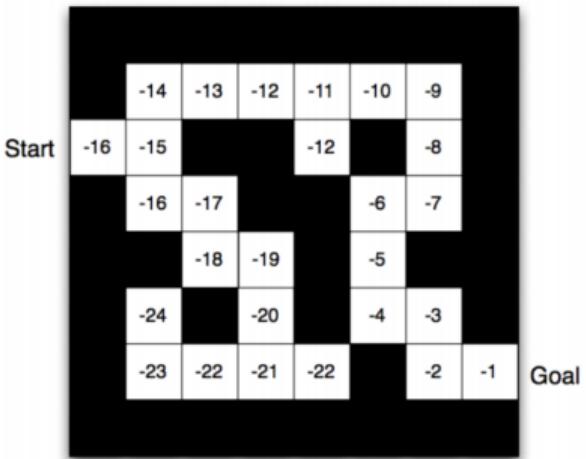
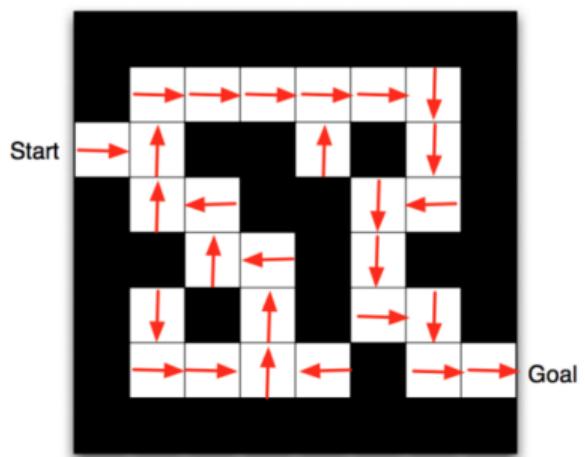
- ▶ How good is a state?

- ▶ How good is a state?
- ▶ The **value function** at state s , is the expected cumulative reward from following the policy from state s :

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, \pi \right]$$

- ▶ Computing the value function is generally impractical, but we can try to approximate (learn) it
- ▶ The benefit is credit assignment: see directly how an action affects future returns rather than wait for rollouts

Value Function



Q-Value Function

- ▶ Can we use a value function to choose actions?

$$\arg \max_a r(s_t, a) + \gamma \mathbb{E}_{p(s_{t+1}|s_t, a_t)} [v^\pi(s_{t+1})]$$

Q-Value Function

- ▶ Can we use a value function to choose actions?

$$\arg \max_a r(s_t, a) + \gamma \mathbb{E}_{p(s_{t+1}|s_t, a_t)} [v^\pi(s_{t+1})]$$

- ▶ **Problem:** this requires taking the expectation with respect to the environment's dynamics, which we don't have direct access to!
- ▶ Instead, learn a Q-value function.

Q-Value Function

- ▶ The Q-value function at state s and action a , is the expected cumulative reward from taking action a in state s and then following the policy:

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

- ▶ Relationship:

$$V^\pi(s) = \sum_a \pi(a|s) Q^\pi(s, a)$$

- ▶ Optimal Action:

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

Bellman Equation

- ▶ The Bellman Equation is a recursive formula for the Q-value function:

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)\pi(a'|s')} [Q^\pi(s', a')]$$

- ▶ There are various Bellman equations, and most RL algorithms are based on repeatedly applying one of them.

Bellman Equations:

- ▶ For the value function $V^\pi(s)$:

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^\pi(s')]$$

- ▶ For the Q-value function $Q^\pi(s, a)$:

$$Q^\pi(s, a) = \sum_{s'} P(s'|s, a) \left[R(s, a, s') + \gamma \sum_{a'} \pi(a'|s') Q^\pi(s', a') \right]$$

Optimal Bellman Equation

- ▶ The optimal policy π^* is the one that maximizes the expected discounted reward, and the optimal Q-value function Q^* is the Q-value function for π^*
- ▶ The Optimal Bellman Equation gives a recursive formula for Q^* :

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} \left[\max_{a'} Q^\pi(s_{t+1}, a') | s_t = s, a_t = a \right]$$

- ▶ This system of equations characterizes the optimal Q-value function. So maybe we can approximate Q^* by trying to solve the optimal Bellman equation!
- ▶ **Intuition:** If the optimal state-action values for the next time-step $Q^*(s', a')$ are known, then the optimal strategy is to take the action that maximizes the expected value of $r(s, a) + \gamma Q^\pi(s', a')$

Reinforcement Learning: Q-Learning

- ▶ Q-learning is a **model-free** reinforcement learning algorithm
- ▶ It **learns the value of actions in states**, known as the **Q-value function**
- ▶ The goal is to find the optimal policy that maximizes the expected cumulative reward
- ▶ Q-learning is off-policy, meaning it learns the value of the optimal policy independently
- ▶ It can be used in environments with unknown dynamics, making it suitable for a wide range of problems

Q-Learning (cont.)

- ▶ Let Q be an action-value function which hopefully approximates Q^*
- ▶ Q-learning is an algorithm that repeatedly adjusts Q to minimize the Bellman error
- ▶ Each time we sample consecutive states and actions (s_t, a_t, s_{t+1}, r_t)

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \underbrace{\left[r(s, a) + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right]}_{\text{Bellman Error}}$$

Exploration-Exploitation Tradeoff

- ▶ Notice: Q-learning only learns about the states and actions it visits.
- ▶ Exploration-exploitation tradeoff: the agent should sometimes pick suboptimal actions in order to visit new states and actions.
- ▶ Simple solution: ϵ -greedy policy
 - With probability $1 - \epsilon$, choose the optimal action according to Q
 - With probability ϵ , choose a random action
- ▶ ϵ -greedy policy still widely used today regardless of simplicity

Q-Learning Algorithm

Initialize $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

 Choose A from S using policy derived from Q (e.g., ε -greedy)

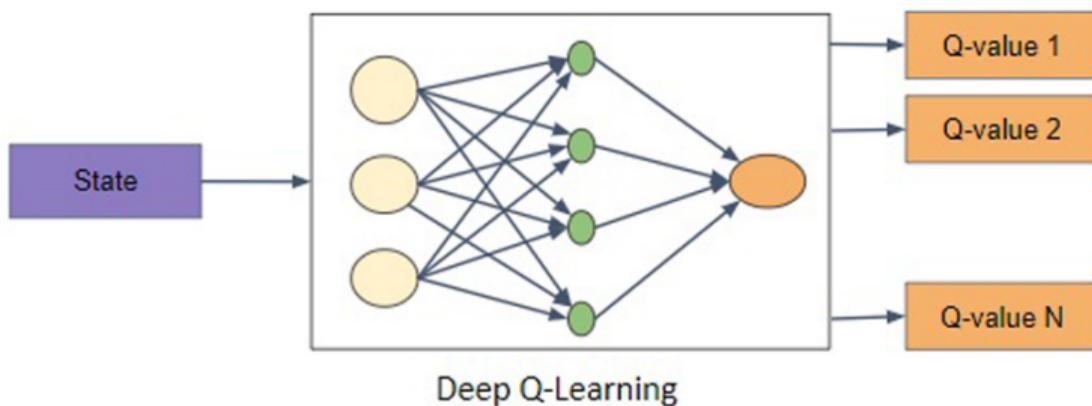
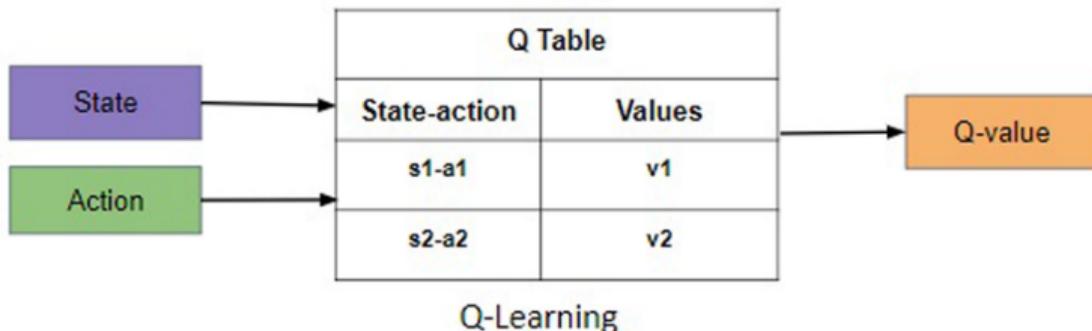
 Take action A , observe R, S'

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

$S \leftarrow S'$;

 until S is terminal

- ▶ Q-learning converges to the optimal action-value function Q^* under certain conditions:
 - **Sufficient exploration:** Every state-action pair is visited infinitely often.
 - **Learning rate:** The learning rate α decays appropriately, i.e., $\alpha \rightarrow 0$ as the number of visits increases, but $\sum_t \alpha_t = \infty$ and $\sum_t \alpha_t^2 < \infty$.
 - **Finite MDP:** The environment is a finite Markov Decision Process.
- ▶ Under these conditions, Q-learning is guaranteed to converge to Q^* with probability 1.



Reinforcement Learning: **Limitations**

Limitations of Reinforcement Learning

- ▶ **Sample inefficiency:** Requires a large number of interactions with the environment to learn effectively.
- ▶ **Stability issues with function approximation:** Training with neural networks can be unstable and sensitive to hyperparameters.
- ▶ **Sparse or delayed rewards:** Learning is difficult when rewards are infrequent or delayed.
- ▶ **Requires careful reward design:** Poorly designed rewards can lead to unintended behaviors.
- ▶ **Poor generalization across tasks:** Policies often do not transfer well to new or slightly different tasks.

Reinforcement Learning: **Summary**

- ▶ Reinforcement Learning (RL) is a powerful paradigm for **sequential decision-making**.
- ▶ Markov Decision Processes (MDPs) provide a **mathematical foundation**.
- ▶ Q-learning enables **model-free control**.
- ▶ Despite its power, RL has **significant limitations**.
- ▶ Research is ongoing to make RL **scalable, robust, and safe**.

Reinforcement Learning: References

References

- [1] Sutton, R. S., & Barto, A. G. (2018). *Reinforcement Learning: An Introduction* (2nd ed.).
- [2] Mnih, V., Kavukcuoglu, K., Silver, D., et al. (2015). Human-level control through deep reinforcement learning. *Nature*.
- [3] Silver, D., Huang, A., Maddison, C. J., et al. (2016). Mastering the game of Go with deep neural networks and tree search. *Nature*.
- [4] Levine, S., Kumar, A., Tucker, G., & Fu, J. (2020). Offline Reinforcement Learning: Tutorial, Review, and Perspectives. *arXiv:2005.01643*.
- [5] David Silver's RL Course:
<http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>

References (cont.)

- [6] Berkeley CS285: Deep RL
<https://rail.eecs.berkeley.edu/deeprlcourse/>

- [7] Emma Brunskill, Stanford CS234: Reinforcement Learning
<https://web.stanford.edu/class/cs234/>

Credits

Dr. Prashant Aparajeya

Computer Vision Scientist — Director(AISimply Ltd)

p.aparajeya@aisimply.uk

This project benefited from external collaboration, and we acknowledge their contribution with gratitude.