# From Classical to Quantum Reinforcement Learning and Its Applications in Quantum Control: A Beginner's Tutorial

Abhijit Sen,[1, *] Sonali Panda,[2, †] Mahima Arya,[1, ‡] Subhajit Patra,[3, §] Zizhan Zheng,[3, ¶] and Denys I. Bondar[1, **]

[1]*Department of Physics and Engineering Physics,*
*Tulane University, New Orleans, LA 70118, USA*
[2]*Department of Physics, Indian Institute of Technology, Dhanbad, India*
[3]*Department of Computer Science, Tulane University, New Orleans, LA 70118, USA*
(Dated: January 14, 2026)

This tutorial is designed to make reinforcement learning (RL) more accessible to undergraduate students by offering clear, example-driven explanations. It focuses on bridging the gap between RL theory and practical coding applications, addressing common challenges that students face when transitioning from conceptual understanding to implementation. Through hands-on examples and approachable explanations, the tutorial aims to equip students with the foundational skills needed to confidently apply RL techniques in real-world scenarios. *Code Availability*: https://github.com/asen009/Reinforcement_Tutorial_Undergraduates

## I. INTRODUCTION

Reinforcement Learning (RL) is a branch of artificial intelligence that trains an agent to maximize desirable outcomes through interaction with its environment. RL has seen rapid advancements in recent years, leading to breakthroughs in various fields. In the field of games [1, 2], autonomous navigation [3, 4], robotics [5] and Large Language Models [6], significant development was possible due to RL. In the near future, AI models are expected to increasingly integrate supervised and unsupervised learning algorithms with RL, creating hybrid approaches that leverage the strengths of each method. It is advised to refer to the following resource as a prerequisite [7].

The study of RL is crucial for students due to its growing applications in various fields [8]. However, one notable shortcoming of the field is the significant time investment required to master its concepts and techniques. Furthermore, it is equally important that the gap between theoretical understanding and code implementation of reinforcement learning is taught side by side so that students understand the actual implementation of what they learn. Many tutorials and journal articles on reinforcement learning provide useful insights[9–11], but they are often lengthy, mathematically heavy, or scattered across multiple examples, making it hard for readers to build a clear understanding. This tutorial takes a different approach by using a single, simple example to explain all the main concepts in a connected and easy-to-follow way. It focuses only on the ideas that are necessary

and sufficient, keeping the content concise but complete. Each method is introduced by addressing what the previous one lacks, helping readers see how the techniques improve step by step. The tutorial also includes clear mathematical explanations along with ready-to-use code, making it practical and accessible for anyone looking to learn reinforcement learning in a structured and efficient manner.

Before explaining how RL can be used for quantum control, we will first present RL in simple yet complete terms. Hence, the rest of the paper is organized as follows: First of all in Sec. II, we provide an overview of the fundamental concepts of RL. Then in III, we have explicitly discussed the concept of probability which is an essential prerequisite. The subsequent sections IV and V delve into the essential terminology commonly used in RL. In Sec. VI, we discuss key elements such as policies, reward systems, and actions. Further, we explore policy improvement techniques, while in Sec. VII we continue the discussion on policy improvement strategies. Section VIII covers the Markov Decision Process (MDP) and its associated techniques. Dynamic Programming is introduced in Sec. IX. Section X presents the average policy evaluation strategy via the Monte Carlo Methods. In Sec. XI, we examine the bootstrapping strategy for policy evaluation via temporal difference methods. Section XII and XIII introduces direct policy optimization techniques, including policy gradient and actor-critic algorithms, which are well-suited for continuous action spaces. Then, Sec. XIV and XV outlines how these reinforcement learning methods enable efficient, high-fidelity manipulation of quantum states, implementing quantum control.

Finally, after concluding remarks in Sec. XVII, the appendices provide detailed discussions on the Bellman optimality condition, iterative methods, iterative policy evaluation, value iteration, as well as both first-visit and every-visit monte carlo methods.

* asen1@tulane.edu, abhijit913@gmail.com
† sp.sonalpanda@gmail.com
‡ marya@tulane.edu, aryamahima@gmail.com
§ subh.p9083@gmail.com
¶ zzheng3@tulane.edu
** dbondar@tulane.edu

## II. UNDERSTANDING RL

The term "reinforcement" in RL refers to the process of encouraging or strengthening a behavior or action. For example, in a park (environment) filled with many games, giving a treat (reward) to a dog (agent) for playing specific type of games (goal) is a form of positive reinforcement. This leads the dog to learn and prefer those specific games in the park. Over time, the dog learns to associate with these reward generating games and adjusts its behavior to play only those games that maximize its treats. With this illustration in mind, we are now in a position to define RL. Reinforcement learning is a type of machine learning where an agent interacts with an environment to achieve a goal by taking actions and receiving feedback in the form of rewards or penalties[12]. The agent's goal is to maximize the cumulative reward over time by learning which actions lead to the most favorable outcomes.

Let us delve into more examples. In the field of autonomous driving, a car, acting as an agent, interacts with its road environment and receives rewards or penalties based on its actions. For example, it earns a reward for successfully avoiding a collision, while it faces a penalty for failing to stop at a red light. Over time, the car learns the best driving behaviors/action for different situations and becomes a successful self-driving car. A successful agent is the one who excels at decision making. Thus, RL helps us in producing good decision-making agents by enabling them to learn from their interactions with the environment. Note that the agent (self-driving car) learns through trials and errors. However, there is a caution. If the car is not made to interact with an environment that is complicated enough to mimic real-life traffic or road scenarios, the RL model may not generalize well when deployed in the real world.

Another example can be taken from the field of suppy-chain management. In supply chain management, through RL we can train an agent that can optimize inventory management and demand forecasting. Here, the agent (the management system) interacts with its environment (market demand, inventory levels, production rates, and transportation costs) to perform a suitable action (the number of items to order on a certain day) in order to balance costs, inventory, and service quality.

Note that RL is different from supervised learning in artificial intelligence. The core difference is, in supervised learning, the model learns from a labeled dataset, where the correct input-output pairs are known, whereas, in RL the model (agent) learns by interacting with an environment and receiving feedback in the form of rewards or penalties based on its actions.

## III. ESSENTIAL PRELIMINARIES

To make this tutorial self-contained, we briefly describe the minimal probabilistic background, random variables and expectations.

**Probability Theory**: Probability quantifies the likelihood of an event occurring and is constrained to the interval $[0, 1]$, where 0 represents an impossible event and 1 represents certainty. For a fair coin toss, we have $\mathbb{P}(\text{Head}) = \mathbb{P}(\text{Tail}) = \frac{1}{2}$. Similarly, for a fair six-sided die, each outcome has equal probability $\mathbb{P}(\text{face } i) = \frac{1}{6}$ for $i \in \{1, 2, 3, 4, 5, 6\}$.

When two events $A$ and $B$ are independent, meaning the occurrence of one does not influence the probability of the other, their joint probability follows the multiplication rule:

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \cdot \mathbb{P}(B) \tag{1}$$

Consider a game where Alice must obtain heads on two consecutive coin tosses to win. Since coin tosses are independent events:

$$\mathbb{P}(\text{Head on toss 1 and Head on toss 2}) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$$

This principle extends naturally to multiple independent events and forms the basis for many statistical models.

When events $A$ and $B$ are mutually exclusive (cannot occur simultaneously), their union probability follows the addition rule:

$$\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) \tag{2}$$

In another game, Alice tosses a single fair six–sided die and wins if the outcome is either 1 or 6. The events "roll a 1" and "roll a 6" are mutually exclusive, so by the addition rule,

$$\mathbb{P}(1 \text{ or } 6) = \mathbb{P}(1) + \mathbb{P}(6) = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}.$$

**Conditional Probability**: Conditional probability captures how our assessment of an event's likelihood changes when we acquire additional information. This concept is fundamental to Bayesian statistics, causal inference, and machine learning algorithms. The conditional probability of event $A$ given event $B$ is defined as:

$$\mathbb{P}(A \mid B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \quad \text{when } \mathbb{P}(B) > 0 \tag{3}$$

Note that, in $\mathbb{P}(A \mid B)$, the vertical bar "|" means "given," i.e., the probability of $A$ conditioned on the event $B$.

Let us understand conditional probability through an example. Suppose in the same game of tossing two coins, Alice receives inside information that the first coin always produces heads irrespective of how it is tossed. In such a scenario, Alice's chances of winning are dramatically improved. Let us examine how conditional probability helps us quantify this advantage.

Given Alice's privileged information about the biased nature of the first coin, she knows with certainty that any outcome must begin with a head. This knowledge

fundamentally transforms her probability calculations. The sample space, which originally contained four equally likely outcomes $\{HH, HT, TH, TT\}$, now reduces to just two possible outcomes: $\{HH, HT\}$. Since the second coin remains fair and independent, both remaining outcomes are equally probable.

If Alice's objective is to achieve two heads (both coins showing heads), her probability of success becomes:

$$\mathbb{P}(\text{both heads} \mid \text{first coin always heads}) = \frac{1}{2}.$$

If the first coin is guaranteed to show heads, then $\mathbb{P}(\text{first coin H}) = 1$ and the joint event "both heads" is equivalent to "second coin is heads." Hence

$$\mathbb{P}(\text{both H} \mid \text{first coin always H}) = \frac{\mathbb{P}(\text{both H} \cap \text{first coin H})}{\mathbb{P}(\text{first coin H})}$$
$$= \frac{\mathbb{P}(\text{second coin H})}{1}$$
$$= \frac{1}{2}.$$

**Random variables**: Consider rolling a fair six-sided die. The die has 6 numbers on its faces: $\{1, 2, 3, 4, 5, 6\}$. When we roll it, a single number appears on the top face. This number changes each time we roll the die again. To capture this changing value mathematically, we define a variable $X$ and call it a *random variable* that represents the number showing on top after any single roll.

The random variable $X$ can take any of the six possible values. Since the die is fair, each outcome is equally likely, so the probability of getting any specific number is the same. For example, $\mathbb{P}(X = 6) = \frac{1}{6}$, meaning there's a one-in-six chance of rolling a 6. More generally, we write $\mathbb{P}(X = k) = \frac{1}{6}$ for $k = 1, 2, 3, 4, 5, 6$. This collection of probabilities is called the probability mass function (pmf), which tells us how likely each possible outcome is.

We can also ask questions like "What is the probability of rolling 3 or less?" This leads to the cumulative distribution function (cdf), defined as $F_X(x) = \mathbb{P}(X \leq x)$. For our die example, $F_X(3) = \mathbb{P}(X \leq 3) = \mathbb{P}(X = 1) + \mathbb{P}(X = 2) + \mathbb{P}(X = 3) = \frac{3}{6} = \frac{1}{2}$.

To summarize the "typical" value we expect from rolling the die, we compute the expectation (or mean):

$$\mathbb{E}[X] = \sum_{k=1}^{6} k \cdot \mathbb{P}(X = k) = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + \cdots + 6 \cdot \frac{1}{6} = \frac{21}{6} = 3.5.$$

**Conditional Expectation**: Conditional expectation extends the concept of averaging to scenarios where we possess additional information. When we know that event $B$ has occurred (with $\mathbb{P}(B) > 0$), the conditional expectation of a discrete random variable $X$ given $B$ computes the average using only those outcomes consistent with the observed information:

$$\mathbb{E}[X \mid B] = \sum_x x \, \mathbb{P}(X = x \mid B)$$
$$= \frac{1}{\mathbb{P}(B)} \sum_x x \, \mathbb{P}(X = x, \, B). \tag{4}$$
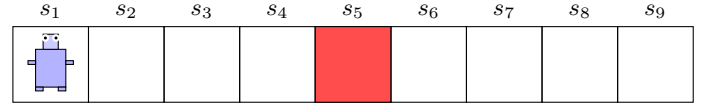


FIG. 1: A 1D grid with a robot in cell $s_1$ and a wall in cell $s_5$.

Returning to our die example where $X \in \{1, \ldots, 6\}$ represents a fair roll, suppose we learn that "the roll is at least 4." This information restricts our attention to outcomes $\{4, 5, 6\}$, each equally likely within this subset. The conditional expectation becomes:

$$\mathbb{E}[X \mid X \geq 4] = \frac{4 \cdot \frac{1}{6} + 5 \cdot \frac{1}{6} + 6 \cdot \frac{1}{6}}{\frac{3}{6}} = \frac{4 + 5 + 6}{3} = 5. \tag{5}$$

With these essential mathematical foundations now established, one now has the theoretical toolkit necessary for understanding reinforcement learning.

## IV. TERMINOLOGY ALERT: PART I

In RL, we need to understand several key terms beyond the basics (agent, action, environment, goals, and rewards). These include policies, transition probabilities, state-value functions, action-value functions, episodes, trajectories, and discount factors. Let's explore each concept using a simple example.

Consider a 1D grid with 9 cells (see Fig. 1). Each cell in the grid is a state, hence there are 9 such states for the robot to be in. Let $\mathcal{S} = (s_1, s_2, .., s_9)$ be the state set. Here, the **agent** is the **robot** and the **1D grid** is the **environment** with which the agent can interact. Further, state $s_5$ is defined to be the terminal state, i.e., the agent cannot cross it.

Now that we have created an agent and the environment, it is time to introduce actions, goal, and reward system to achieve the goal. In the 1D grid environment, the agent in any state is allowed to take any of the two **actions** – $\mathcal{A} = (\texttt{"right"}, \texttt{"left"})$. Note that if the agent takes action $\texttt{"left"}$ by being in state $s_1$ or action $\texttt{"right"}$ in state $s_9$, the state is the same as the current state. The goal of the agent is to reach the terminal state $s_5$ following a policy (defined below). To determine which actions are preferable, we introduce a reward function that assigns a numerical reward or penalty to each action taken and its resulting outcome. Our objective is to encourage the agent to reach a terminal state as quickly as possible. Therefore, we can assign a negative reward (penalty) for actions that move the agent away from the goal or increase the time to reach it, and a positive reward for successfully completing the task. Additionally, we can define a cumulative reward over an entire trajectory (episode), and train the agent to maximize the total expected return–that is, the sum of rewards collected throughout the episode. Hence, the **reward system** is

| **State** | $\mathcal{P}(right)$ | $\mathcal{P}(left)$ |
|-----------|------|------|
| $s_1$ | 1 | 0 |
| $s_2$ | 1 | 0 |
| $s_3$ | 1 | 0 |
| $s_4$ | 1 | 0 |
| $s_5$ | 0 | 0 |
| $s_6$ | 1 | 0 |
| $s_7$ | 1 | 0 |
| $s_8$ | 1 | 0 |
| $s_9$ | 1 | 0 |

TABLE I: Policy $\pi$ mapping for the 1D grid example.

chosen such that the agent gets a reward of -1 for each step it takes toward any regular state, but receives +5 when it reaches the terminal (goal) state.

**Policy**: Now let us design a **policy** for the agent to achieve this goal. We will start with the formal mathematical definition, then illustrate it with a concrete example.

A policy $\pi(a|s)$ is a mapping from states to probabilities over actions,

$$\pi : \mathcal{S} \to \mathcal{P}(\mathcal{A}) \tag{6}$$

where $\mathcal{S}$ is the set of all possible states, $\mathcal{A}$ is the set of all possible actions, and $\mathcal{P}(\mathcal{A})$ is the set of probability distributions over $\mathcal{A}$. Thus, for any state $s$, the policy $\pi(a|s)$ gives the probability of choosing action $a$ in that state. The policy essentially defines a probability distribution over actions for each state, satisfying

$$\sum_{a \in \mathcal{A}} \pi(a|s) = 1 \quad \text{for all } s \in \mathcal{S}. \tag{7}$$

We start with a deterministic policy specified by Table I. For the toy example of Fig. 1, it should be obvious that such a policy is not optimal. However, we will show below how to improve it iteratively.

The information in Table I can be simply written in the form of a Python dictionary as follows,

```python
policy = {
    "s1": {"right": 1, "left": 0},
    "s2": {"right": 1, "left": 0},
    "s3": {"right": 1, "left": 0},
    "s4": {"right": 1, "left": 0},
    "s5": {},   # Terminal state
    "s6": {"right": 1, "left": 0},
    "s7": {"right": 1, "left": 0},
    "s8": {"right": 1, "left": 0},
    "s9": {"right": 1, "left": 0}
}
```

Note that the policy was implemented in Python using the built-in dictionary [13] data type. A dictionary stores data as key-value pairs, where each key maps to its corresponding value. For instance, `python"s1"` is a key that maps to a nested dictionary `python"right": 1, "left": 0`. The key advantage of dictionaries is their efficient value retrieval: given a key, the corresponding value can be

extracted in constant time, regardless of the dictionary's size! For example, the command `pythonpolicy["s1"]` returns the associated value with $O(1)$ time complexity. This efficiency is achieved through hash-based storage, which does not preserve insertion order [14].

For the policy in Table I, $\pi(right|s_2) = 1$ and $\pi(left|s_2) = 0$ and therefore we have

$$\sum_{a \in \mathcal{A}} \pi(a|s_2) = \pi(right|s_2) + \pi(left|s_2) = 1. \tag{8}$$

Note that the policy serves as a guiding framework for the agent, indicating which action to take when it is in a specific state. In our case, the policy (see Table I) is deterministic because it allows the agent to take the action `"right"` for any state it is in. However, at this point, a natural question arises: What is the probability that this action will lead the agent to end up in state $s_2$? In RL, uncertainty is inherent in the agent's interactions with its environment. For instance, when the agent is in state $s_1$ and decides to take the action `"right"`, we cannot definitively conclude that it will end up in state $s_2$. This uncertainty arises because we have not previously specified the potential outcomes associated with the actions available to the agent in each state. This question leads us to the concept of **state-transition probabilities**, denoted as $p(s'|s,a)$.

The state-transition probability $p(s'|s,a)$ is a three-argument function:

$$\mathcal{P} : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \to [0,1]. \tag{9}$$

This represents the probability that the agent transitions to state $+ s'$ given current state $s$ and action $a$. In other words, for any state-action pair $(s,a)$, $p(s'|s,a)$ gives the probability of moving to the next state $s'$. As a probability distribution, it must satisfy

$$\sum_{s' \in \mathcal{S}} p(s'|s,a) = 1. \tag{10}$$

To simplify our analysis, we use deterministic state transitions where probabilities are either 0 or 1. When an agent takes an action from a given state, it transitions to a single predetermined next state with probability 1, while all other transitions have probability 0.

Consider the example of state transition probabilities in Table II. Due to the deterministic nature, the transition probability can be represented by a Python dictionary.

```python
transitions = {
    "s1": {"right": "s2", "left": "s1"},
    "s2": {"right": "s3", "left": "s1"},
    "s3": {"right": "s4", "left": "s2"},
    "s4": {"right": "s5", "left": "s3"},
    "s5": {},   # Terminal state
    "s6": {"right": "s7", "left": "s5"},
    "s7": {"right": "s8", "left": "s6"},
    "s8": {"right": "s9", "left": "s7"},
    "s9": {"right": "s9", "left": "s8"}
}
```

| State | Action | Next State ($s'$) | Probability |
|-------|--------|------------------|-------------|
| $s_1$ | right | $s_2$ | 1 |
| $s_1$ | left | $s_1$ | 1 |
| $s_2$ | right | $s_3$ | 1 |
| $s_2$ | left | $s_1$ | 1 |
| $s_3$ | right | $s_4$ | 1 |
| $s_3$ | left | $s_2$ | 1 |
| $s_4$ | right | $s_5$ | 1 |
| $s_4$ | left | $s_3$ | 1 |
| $s_5$ | right | $s_5$ | 1 |
| $s_5$ | left | $s_5$ | 1 |
| $s_6$ | right | $s_7$ | 1 |
| $s_6$ | left | $s_5$ | 1 |
| $s_7$ | right | $s_8$ | 1 |
| $s_7$ | left | $s_6$ | 1 |
| $s_8$ | right | $s_9$ | 1 |
| $s_8$ | left | $s_7$ | 1 |
| $s_9$ | right | $s_9$ | 1 |
| $s_9$ | left | $s_8$ | 1 |

TABLE II: Intuitive State Transition Probabilities $P(s'|s, a)$ indicating next states and their probabilities.
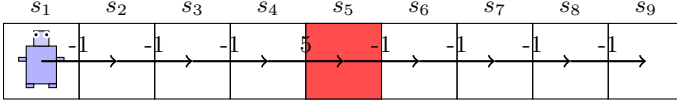


FIG. 2: An improved version of Figure 1 with action and rewards explicitly shown.

It is natural to wonder what happens when we have non-deterministic transitions and what might cause that. For example, if $p(s_2|s_1, right) = 0.7$ and $p(s_4|s_1, right) = 0.3$, then taking action "right" from state $s_1$ leads to state $s_2$ with 70% probability and to state $s_4$ with 30% probability.

This uncertainty might come from a "mysterious force" in the environment, which adds randomness to what happens next. Even though the agent can choose its actions based on its current state and policy, the outcome is not guaranteed. So, while the chance of ending up in state $s_4$ is lower at 30%, it is still possible, showing that decision-making can be complicated in uncertain situations. In real-life scenarios, RL environments are often stochastic due to various uncertainties and complexities inherent in the real world. For example, an RL environment where human decisions are involved and given the unpredictability of human behavior and decision, the environment can itself become stochastic. In reinforcement learning, the agent interacts with the environment and receives rewards, which are numerical values indicating how good or bad a taken action was. These rewards guide learning by encouraging the agent to choose actions that maximize long-term benefit.

Before going further, we graphically represent the introduced concepts in Fig. 2.

**Episode and Trajectory**: Let us assume that our agent in the 1D grid world follows the policy $\pi$. The goal is to reach the terminal state $s_5$. If the agent begins from the current state $s_1$ as shown in Figure 2 and follows the policy $\pi$, it reaches the terminal state $s_5$ in a finite number of steps. In this context, an episode begins when the agent starts in state $s_1$. The episode progresses as the agent follows the actions dictated by policy $\pi$ and moves through various states in the grid. The episode continues until the agent successfully reaches the terminal state $s_5$. Throughout the episode, the agent collects rewards.

The **trajectory**, on the other hand, demonstrates how the agent navigates through the grid world with the decision it makes at each state according to policy $\pi$ until it reaches the terminal state $s_5$ where the episode ends. A **trajectory** can be described as a set that contains tuples of the form $(s_t, a_t, r_{t+1})$, where $s_t$ represents the current state of the agent at time $t$, $a_t$ denotes the action taken by the agent in state $s_t$, $r_{t+1}$ signifies the reward received after taking action $a_t$. In our case, for the episode where the agent begins in the current state $s_1$ and reaches the terminal state $s_5$, the trajectory is (in Python list and tuples):

```
Trajectory = [
    ("s1", "right", -1),
    ("s2", "right", -1),
    ("s3", "right", -1),
    ("s4", "right", 5)
]
```

To connect this with the mathematical notation, we align the Python labels with the time-indexed states as follows:

$$\begin{aligned} \text{"s1"} &\longleftrightarrow s_0 \quad \text{(initial state at } t=0) \\ \text{"s2"} &\longleftrightarrow s_1 \\ \text{"s3"} &\longleftrightarrow s_2 \\ \text{"s4"} &\longleftrightarrow s_3 \end{aligned}$$

Thus, a trajectory is

$$[(s_0, a_0, r_1), (s_1, a_1, r_2), (s_2, a_2, r_3), (s_3, a_3, r_4)]$$

**State Value Function**: The state value function of a state $s$ under a policy $\pi$, denoted $v_\pi(s)$, is the expected return when starting in $s$ and following policy $\pi$ thereafter [15]. Let us understand the definition in our context. Given the agent is in state $s_1$ and follows policy $\pi$, the return is the sum of rewards we see in the trajectory above. Note that **reward** refers to the immediate feedback received after each action, while **return** $G_t$ represents the cumulative return from time $t$ onward, defined as the sum of future rewards received by the agent (or sum of future rewards starting from current state state until the end of the episode in an episodic task)

$$G_t = r_{t+1} + r_{t+2} + ... + r_T. \tag{11}$$

where $T$ is the final time step that leads the agent to the terminal state. In our case,

$$\begin{aligned} \text{Return} &= \text{Sum of future rewards} \\ &= (-1) + (-1) + (-1) + (+5) = +2. \end{aligned}$$

Since the trajectory is deterministic (there is no variability in the trajectory across different episodes starting from the current state, say $s_1$), in every run, the sum of reward/return is the same. Therefore, the value of a state in this context can be directly computed as the total return from that state to the end of the episode, and the expectation becomes unnecessary. Thus, the value of a function for the deterministic case as ours can be mathematically defined as

$$v_\pi(s) = G_t|_{S_t=s} \qquad (12)$$

where $S_t = s$ indicates the initial state of the agent/state from which the agent begins following the policy.

Let's examine which state is more beneficial for the agent to be in. A quick calculation reveals that state $s_4$ has a value of $+5$, which is higher than any other state. This value is expected to be elevated since it is closer to the terminal state.

Expectation is essential when the environment is stochastic, which means there is randomness in the transitions or rewards, or both. This randomness can lead to various possible outcomes for the agent's trajectory, even if it starts from the same state and follows the same policy every time. In such a case, the value of a state can be defined mathematically (take $\gamma = 1$) as:

$$v_\pi(s) = \mathbb{E}_\pi \left[ \sum_{k=0}^{T-t-1} r_{t+k+1} \mid S_t = s \right] \qquad (13)$$

The state-value function under a policy $\pi$, denoted $v_\pi(s)$, represents the expected return when starting in state $s$ and following policy $\pi$ thereafter. It is formally defined as:

$$v_\pi(s) \doteq \mathbb{E}_\pi \left[ G_t \mid S_t = s \right] = \mathbb{E}_\pi \left[ \sum_{k=0}^{T-t-1} \gamma^k r_{t+k+1} \,\middle|\, S_t = s \right], \qquad (14)$$

where: $\gamma \in [0,1]$ is the discount factor, $r_{t+k+1}$ is the reward received at time step $t + k + 1$, $G_t = \sum_{k=0}^{T-t-1} \gamma^k r_{t+k+1}$ is the return from time $t$, and $\mathbb{E}_\pi[\cdot]$ denotes the average / expectation over all possible trajectories generated by following policy $\pi$, i.e., it accounts for both the stochasticity of the environment and the action selection governed by $\pi$.

**Discount Factor ($\gamma$):** Up until now, we discussed agent-environment interaction that can be conceptualized in the form of episodes. In such cases, there is a clear starting point and an ending point (when the agent hits the terminal state). However, there are cases when agent-environment interaction can go on continuously without a well-defined ending point. In cases like these the idea of episodes does not apply. An example of this ongoing, non-episodic interaction is temperature control in an industrial environment, where the system constantly adjusts heating or cooling levels to maintain optimal conditions. Here, there is no natural "end" to the process—the controller operates continuously without a reset, adjusting based on a feedback loop that runs indefinitely.

In cases of continuous interaction between the agent and environment, defining the return $G_t$ can be challenging because it may grow indefinitely. To keep the return manageable, we use a discount factor, which ensures that $G_t$ stays finite by reducing the impact of future rewards over time. In such cases, we have,

$$G_t \doteq r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}, \quad (15)$$

where an extra $\gamma$ factor is introduced such that for $\gamma < 1$, the above sum stays finite.

## V. TERMINOLOGY ALERT: PART II

In the above section, we explained basic terminology a student encounters in studying RL. In this section, we will explain other terminologies, such as state-reward transition probability and action-value functions.

**State-Reward Transition Probability**: In the previous section, we understood state-transition probability $p(s'|s,a)$ which gives the probability for an agent to reach state $s'$ given it is in state $s$ and takes action $a$ being in the current state. Similarly, state-reward transition probability $p(s', r|s, a)$ is the probability for the agent to reach state $s'$ and simultaneously achieve reward $r$ given it is in state $s$ and takes action $a$ being in the current state. The relation between both quantities is as follows:

$$\sum_r p(s', r|s, a) = p(s'|s, a) \qquad (16)$$

In our example of 1D grid, $p(s_2, 5|s_1, right) = 0$ but $p(s_2|s_1, right) = 1$. This is because, in the former case, a transition from $s_1$ to $s_2$ does take place when the action taken is right but the reward is not 5 instead it is -1. In the latter case, it is straightforward to see that it is indeed true. Note that the function $p(s'|s, a)$ shows the likelihood of reaching a new state but offers no information about rewards, which the agent seeks to maximize. Without reward information, the agent may end up in desirable states without knowing which actions provide the most benefit. In contrast, $p(s', r|s, a)$ gives the agent a clearer picture by considering both next states and their associated rewards, enabling better decision-making.

The environment's uncertainty is encapsulated in the state-reward transition probability $p(s', r \mid s, a)$. In a deterministic setting, this probability is either 0 or 1, indicating a certain outcome. In a non-deterministic environment, it includes probability values and becomes non-trivial.

**Action-Value Function**: In the previous section, we explained value function of a state. The action-value function, denoted $q_\pi(s, a)$, is the expected return when
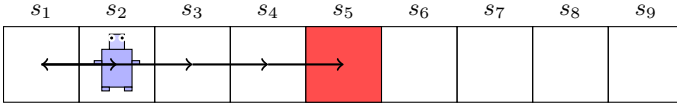
FIG. 3: Another version of Fig. 1 with robot in cell 2 taking first action to left and then following policy $\pi$



FIG. 4: Updated version of Fig. 2 with improved policy

starting in state $s$, taking a specific action $a$, and then following the policy $\pi$ from the resulting state. Mathematically, the action-value function is

$$q_\pi(s,a) = \mathbb{E}_\pi\left[G_t \mid S_t = s, A_t = a\right] \tag{17}$$

$$= \mathbb{E}_\pi\left[\sum_{k=0}^{T-t-1} \gamma^k r_{t+k+1} \,\middle|\, S_t = s, A_t = a\right]. \tag{18}$$

Let us understand this through an example. Note that in the deterministic case, expectation is not necessary (as explained before). In our 1D grid example, assume that the agent is allowed to take action `"left"`. Now if the agent being is state $s_2$ takes the first action `"left"` and then follow the policy (see Fig. 3), then in such case the trajectory is

```
Trajectory = [
    ("s2", "left", -1),
    ("s1", "right", -1),
    ("s2", "right", -1),
    ("s3", "right", -1),
    ("s4", "right", -1),
    ("s5", "right", 5)
]
```

Note that since there is a natural temporal order in a trajectory, we implemented pythonTrajectory as a Python list [16] of tuples [17]. Lists allow efficient extraction of elements (i.e., within the $O(1)$ time) by their ordinal index. However, Python indexing starts from 0. For example, pythonTrajectory[0] returns python("s1", "right", -1), while pythonTrajectory[3] returns python("s4", "right", 5).

From the trajectory we have a return $G_t = 1$. Therefore,

$$q_\pi(s_1, \text{left}) = 1,$$

because the only observed return following action left in state $s_1$ is 1.

There exists a direct relationship between the two value functions. Mathematically, they are related as follows:

$$v_\pi(s) = \sum_a \pi(a|s) q_\pi(s,a) \tag{19}$$

A crucial difference between state-value function and action-value function is that the former is useful for evaluating the value of states, regardless of the action, as it averages over all possible actions according to the policy $\pi$ and the latter is helpful when an agent needs to decide between different actions in a given state since it assesses each action's potential return individually [18].
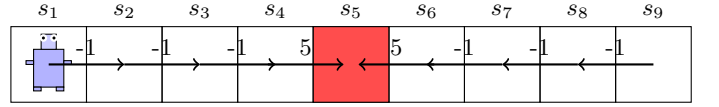
## VI. OBJECTIVE OF RL PART I: SETTING THE STAGE

The primary goal of reinforcement learning (RL) is to enable an agent to learn effective decision-making. This is accomplished by having the agent interact with the environment and use the reward system to refine its policies until it reaches an optimal policy, which ultimately enhances its decision-making abilities.

Let us illustrate policy improvement with an example based on our 1D grid scenario. The original policy allows the agent to reach the terminal state $s_5$ when the agent's current state is any position to the left of $s_5$. However, if the agent finds itself in any state to the right of the terminal state, the current policy does not provide a way to reach the terminal state. Consequently, this policy is not effective and requires improvement. A straightforward enhancement to the policy would be to ensure that when the agent is in any state to the right of the terminal state, the policy allows only the `"left"` action, instead of permitting the `"right"` action. Thus the improved policy (see also Fig. 4) is

```
improved_policy = {
    "s1": {'right': 1, 'left': 0},
    "s2": {'right': 1, 'left': 0},
    "s3": {'right': 1, 'left': 0},
    "s4": {'right': 1, 'left': 0},
    "s5": {},    # Terminal state
    "s6": {'right': 0, 'left': 1},
    "s7": {'right': 0, 'left': 1},
    "s8": {'right': 0, 'left': 1},
    "s9": {'right': 0, 'left': 1}
}
```

At this point, a natural question to ask is: How can we improve the policy? The above example was trivial; hence, quick intuition helped us figure out what might be a good policy. However, in more complicated cases, we need algorithms to assist us in the policy improvement process. Before we move on to algorithms, let us understand what fundamentally dictates the policy improvement process.

The value functions defined previously can be used to dictate what constitutes a better policy. A policy $\pi$ is at least as good as policy $\pi'$ (denoted $\pi \succeq \pi'$) if and only if [15]:

$$v_\pi(s) \geq v_{\pi'}(s), \quad \text{for all } s \in \mathcal{S}. \tag{20}$$

As the agent explores the environment, it improves the values of states (by exploring and updating the values of states based on received rewards) in order to refine the

policy. As the iterative process continues, the agent eventually reaches an optimal policy, where it has maximized its expected reward from each state. Thus, an optimal policy is the policy that maximizes the value of all states. Mathematically, for an optimal policy $\pi = *$, we have

$$v_*(s) = \max_{\pi} v_{\pi}(s), \quad \text{for all } s \in \mathcal{S}, \tag{21}$$

where $v_*(s)$ is the optimal state-value function that one achieves when the policy is optimal. In our case, if the improved policy described above is an optimal policy, then the value of all states $s$ is the optimal value for those states. It is easy to list down the optimal value of the states in this case (note the value of a terminal state is always zero as there is no future reward to collect from there):

```python
state_values = {
    "v(s1)": 2,
    "v(s2)": 3,
    "v(s3)": 4,
    "v(s4)": 5,
    "v(s5)": 0,    # Terminal state
    "v(s6)": 5,
    "v(s7)": 4,
    "v(s8)": 3,
    "v(s9)": 2
}
```

It's important to recognize that optimal policies also have the same optimal action-value function, denoted $q_*$, and defined as [15]:

$$q_*(s, a) \doteq \max_{\pi} q_{\pi}(s, a), \tag{22}$$

for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$, where, $\mathcal{A}(s)$ is the set of actions available in state $s$.

## VII. OBJECTIVE OF RL PART II: POLICY EVALUATION AND POLICY IMPROVEMENT

When an agent lacks complete information about its environment, it learns by interacting with it, often using a trial-and-error approach to improve its policy. In these situations, the agent may encounter multiple possible episodes that can start from the same state, say $s_2$. To find the value of that state, the agent needs to look at the outcomes of all these episodes. By adding up the values from each episode and averaging them, the agent can better understand the state $s_2$ value. In the 1D grid example, it is possible to have many episodes starting from a given state for non-trivial state-transition probabilities. For example, a non-trivial transition with probability (in Python dictionary) is as follows:

```python
transitions = {
  "s1": {
      "right": [("s2", 0.7), ("s1", 0.3)
          ],
      "left": [("s1", 1.0)]
```

```python
  },
  "s2": {
      "right": [("s3", 0.8), ("s1", 0.2)
          ],
      "left": [("s1", 1.0)]
  },
  "s3": {
      "right": [("s4", 0.9), ("s2", 0.1)
          ],
      "left": [("s2", 1.0)]
  },
  "s4": {
      "right": [("s5", 1.0)],
      "left": [("s3", 1.0)]
  },
  "s5": {}, #Terminal state
  "s6": {
      "right": [("s7", 0.6), ("s5", 0.4)
          ],
      "left": [("s5", 1.0)]
  },
  "s7": {
      "right": [("s8", 0.7), ("s6", 0.3)
          ],
      "left": [("s6", 1.0)]
  },
  "s8": {
      "right": [("s9", 0.8), ("s7", 0.2)
          ],
      "left": [("s7", 1.0)]
  },
  "s9": {
      "right": [("s9", 1.0)],
      "left": [("s8", 1.0)]
  }
```

This clearly states that the 1D grid environment is now more stochastic in nature. For example, given an agent is state $s_2$ takes action `"right"`, the agent has some probability (20%) that it can end up in state $s_1$. With this uncertainty, every episode (and its return) starting from a given state will be different.

We will now focus on an environment where complete information is available, specifically regarding the transition probabilities $p(s', r \mid s, a)$. This means we have a thorough understanding of how the environment responds to the agent's actions, including the resulting states and associated rewards. In this case, we do not need to resort to trial and error to determine state values. Instead, we can utilize what is known as the Bellman equation to find the value of each state as shown in Appendix B. The Bellman equation is expressed as follows:

$$v_{\pi}(s) = \sum_{a} \pi(a \mid s) \sum_{\substack{s' \in \mathcal{S} \setminus \{s\} \\ r \in \mathcal{R}}} p(s', r \mid s, a) \left[ r + \gamma v_{\pi}(s') \right]$$

$$\tag{23}$$

for all $s \in \mathcal{S}$, and $s' \in \mathcal{S}$ such that $s' \neq s$. This can be written in a mathematical language as $\forall s \in \mathcal{S}, \quad s' \in \mathcal{S} \setminus \{s\}$.

To understand the Bellman equation Eq. (23) and its significance, consider an environment where the transi-

tion probabilities $p(s', r \mid s, a)$ are known. Given a policy and assuming that the values of all states $s'$ are known (i.e., $v_\pi(s')$ are known for $s' \in \mathcal{S} \setminus \{s\}$), we can solve the Bellman equation to determine the value of a particular state $s$.

Let us illustrate this through an example. Consider a $2 \times 2$ grid with state set $\mathcal{S} = (s_1, s_2, s_3, s_4)$. Further $\mathcal{S} \setminus \{s_1\} = (s_2, s_3, s_4)$, $\mathcal{S} \setminus \{s_2\} = (s_1, s_3, s_4)$, $\mathcal{S} \setminus \{s_3\} = (s_1, s_2, s_4)$ and $\mathcal{S} \setminus \{s_4\} = (s_1, s_2, s_3)$ . Let $s_4$ be the terminal state. The reward and action sets are $\mathcal{R} = (-1, 3)$ and $\mathcal{A} = (\texttt{"right"}, \texttt{"left"}, \texttt{"up"}, \texttt{"down"})$ respectively. Further, the policy and transitions are as follows:

```
policy = {
    "s1": {'right': 0.5, 'down': 0.5},
    "s2": {'left': 0.5, 'down': 0.5},
    "s3": {'right': 0.5, 'up': 0.5},
    "s4": {} #Terminal state
}

  transitions = {
    "s1": {"right": "s2", "down": "s3"},
    "s2": {"left": "s1", "down": "s4"},
    "s3": {"right": "s4", "up": "s1"},
    "s4": {} #Terminal state
}
```

Intuitively, it is clear that the policy is not an optimal policy and needs improvement. However, in order to solve the Bellman equation, we will use this policy (with discount factor $\gamma = 1$).

We will solve the Bellman equation analytically (note that the Bellman equation in practice is solved numerically). Before solving the equation, let us note the following: In Eq. (23), for the above example, we have $\pi(right|s_1) = \pi(down|s_1) = 0.5$ and $\pi(left|s_1) = \pi(up|s_1) = 0$. Similarly, for other states, these quantities can be obtained by looking at the policy given above in Python dictionary format. Likewise, for the transition probabilities, note that for state $s_1$, $p(s_2, -1|s_1, right) = p(s_3, -1|s_1, \text{down}) = 1$, while all other transitions from $s_1$ have zero probability. Transition probabilities for other states can be similarly obtained by referring to the dictionary format provided above.

$$
\begin{aligned}
v_\pi(s_1) &= \sum_{a \in \mathcal{A}} \pi(a \mid s_1) \sum_{\substack{s' \in \mathcal{S} \setminus \{s_1\} \\ r \in \mathcal{R}}} p(s', r \mid s_1, a) \left[ r + v_\pi(s') \right] \\
&= \sum_{a \in \mathcal{A}} \pi(a \mid s_1) \Big\{ p(s_2, -1 \mid s_1, a) \left[ -1 + v_\pi(s_2) \right] \\
&\quad + p(s_3, -1 \mid s_1, a) \left[ -1 + v_\pi(s_3) \right] \Big\} \\
&= 0.5 \Big\{ 1 \left[ -1 + v_\pi(s_2) \right] + 1 \left[ -1 + v_\pi(s_3) \right] \Big\} \\
&= 0.5 \left[ -1 + v_\pi(s_2) \right] + 0.5 \left[ -1 + v_\pi(s_3) \right] \\
&= -1 + 0.5 v_\pi(s_2) + 0.5 v_\pi(s_3)
\end{aligned}
\tag{24}
$$

Similarly, for state $s_3$, we have

$$
\begin{aligned}
v_\pi(s_3) &= \sum_{a \in \mathcal{A}} \pi(a \mid s_3) \sum_{\substack{s' \in \mathcal{S} \setminus \{s_3\} \\ r \in \mathcal{R}}} p(s', r \mid s_3, a) \left[ r + v_\pi(s') \right] \\
&= \pi(a = right \mid s_3) \left[ p(s_4, 3 \mid s_3, right) (3 + v_\pi(s_4)) \right] \\
&\quad + \pi(a = \text{up} \mid s_3) \left[ p(s_1, -1 \mid s_3, \text{up}) (-1 + v_\pi(s_1)) \right] \\
&= 0.5 \left[ 3 + v_\pi(s_4) \right] + 0.5 \left[ -1 + v_\pi(s_1) \right] \\
&= 1.5 + 0.5 v_\pi(s_1) - 0.5 + 0.5 v_\pi(s_4) \\
&= 1 + 0.5 v_\pi(s_1) + 0.5 v_\pi(s_4)
\end{aligned}
\tag{25}
$$

Thus, Bellman equation gives us a set of four equations as follows:

$$
\begin{aligned}
v_\pi(s_1) &= -1 + 0.5 v_\pi(s_2) + 0.5 v_\pi(s_3) \\
v_\pi(s_2) &= 1 + 0.5 v_\pi(s_1) + 0.5 v_\pi(s_4) \\
v_\pi(s_3) &= 1 + 0.5 v_\pi(s_1) + 0.5 v_\pi(s_4) \\
v_\pi(s_4) &= 0
\end{aligned}
\tag{26}
$$

Solving the latter, we find

$$
\begin{aligned}
v_\pi(s_1) &= 0 & v_\pi(s_2) &= 1 \\
v_\pi(s_3) &= 1 & v_\pi(s_4) &= 0
\end{aligned}
\tag{27}
$$

The values of states $s_2$ and $s_3$ are identical, as both states transition to the terminal state $s_4$ in a single step, resulting in the same reward. This equivalence highlights the symmetry in their respective actions and outcomes within the defined environment.

**Policy Evaluation**: Let us return to the Bellman equation (23). We find that the Bellman equations are instrumental in determining the value of states for a specific policy, provided the dynamics of the environment are known.

Policy evaluation refers to the process of calculating the value function for a given policy using the Bellman equations. Policy evaluation is a crucial step in reinforcement learning that involves calculating the value of each state under a specific policy. By performing policy evaluation, we gain insights into how good the current policy is and whether improvements can be made.

**Policy Improvement**: Policy improvement refers to the criteria used to enhance a given policy in reinforcement learning. The aim is to find a policy that yields higher expected returns from each state. An improved policy is one where the values of the states, denoted as $v_\pi(s)$, are greater than those from the previous policy. We will elaborate more on this in Sec. IX.

## VIII. MARKOV CHAIN AND MARKOV DECISION PROCESS (MDP)

In this section, we examine an important assumption about the environment that leads to transition probabilities to be of the form $p(s', r \mid s_1, a)$. We assume that the probability of transitioning to the next state depends

only on the current state and action, not on the agent's previous history. This "memoryless" property is called the Markov property, making our environment a Markovian environment.

In other words, given a trajectory of the form

$$\text{Trajectory} = [(s_{t-2}, a_{t-2}, r_{t-1}), (s_{t-1}, a_{t-1}, r_t), (s_t, a_t, r_{t+1}),$$
$$(s_{t+1}, a_{t+1}, r_{t+2}), (s_{t+2}, a_{t+2}, r_{t+3}), (s_{t+3}, a_{t+3}, r_{t+4})].$$

the transition of the agent from state $s_t$ to $s_{t+1}$ does not depend on the past history. Mathematically,

$$p(s_{t+1} \mid s_t, s_{t-1}, s_{t-2}) = p(s_{t+1} \mid s_t) \qquad (28)$$

RL revolves around decision-making. When we operate within an environment that adheres to the Markov property, we are essentially engaging in a Markov Decision Process (MDP). An MDP is defined by a tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$, where:

- $\mathcal{S}$ represents the set of states in the environment,

- $\mathcal{A}$ is the set of possible actions,

- $P$ transition probability matrix,

- $\mathcal{R}$ is the reward function, and

- $\gamma$ is the discount factor, which balances the importance of immediate versus future rewards.


## IX. DYNAMIC PROGRAMMING (DP)

Dynamic programming (DP) is a mathematical optimization method for finding optimal solutions to Markov Decision Processes (MDPs) [15]. Using recursive equations, DP efficiently improves policies through techniques like policy iteration and value iteration.

In Sec. (VII), we discussed about policy evaluation and policy improvement. Let us delve a bit deeper into how policy improvement is implemented and then we will understand policy iteration. Policy improvement aims to find a new, better policy $\pi'$ than the current policy $\pi$. The new policy $\pi'$ is considered better if the value function for all states under $\pi'$ is at least as high as under $\pi$. Mathematically, this requires that $v_{\pi'}(s) \geq v_\pi(s)$ for all $s$ in $\mathcal{S}$, with strict improvement in at least one state. This means we should seek actions for the agent at each state $s$ that could improve the value function. The process of seeking the best action for a given state is guided by action-value function $q_\pi(s, a)$. The steps to follow are:

1. Given the agent is in state $s$, choose a new action $a$ that the agent takes in current state and then follows policy $\pi$. We can, in this case, calculate state-action value $q_\pi(s, a)$.

2. If

$$q_\pi(s, a) > v_\pi(s),$$

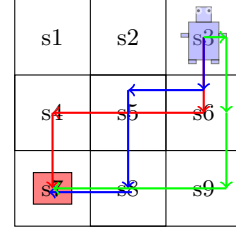

FIG. 5: Robot's trajectories from cell s3 to the target cell s7 with different paths represented by colored arrows.

it suggests that taking action $a$ and then following the current policy is better than just following the policy $\pi$ all the time from that state.

3. Update the policy to choose $a$ in $s$.

Let us revisit our one-dimensional example illustrated in Figure 2, where the policy and transitions are detailed in Sec. IV. For state $s_6$, the existing policy prescribes the action `"right"`, resulting in a state-value of $-3$. However, if we were to assign a new action `"left"` for state $s_6$, we could transition to the terminal state $s_5$ (the terminal state) and receive a action-value of $+3$. In some case, the new action may not lead to terminal state directly and, in such case, the agent is supposed to follow the old policy until the episode ends. This change indicates that by adopting the new action, we significantly improve the action-value of state $s_6$. Consequently, we can update the policy to include this new action for $s_6$ as part of our policy improvement process.

The final step is policy iteration. Once we improve a policy $\pi$ to obtain a better policy $\pi'$ using its value function $v_\pi$, we can then find the value of $\pi'$ as $v_{\pi'}$ and improve it further to get a new policy $\pi''$. This creates a chain of increasingly better policies and value functions [15]:

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \cdots \xrightarrow{I} \pi_* \xrightarrow{E} v_*,$$

where $\xrightarrow{E}$ represents policy evaluation and $\xrightarrow{I}$ represents policy improvement. Each step results in a better policy until we reach the optimal policy.

A drawback of policy iteration is that each iteration requires a complete policy evaluation, which may involve multiple sweeps through the entire state space to achieve convergence. This iterative process can be computationally expensive, especially in environments with a large number of states. In appendices C and D, we introduce value iteration, an alternative approach that integrates policy evaluation and improvement into a single update step, often resulting in faster convergence.

## X. MONTE CARLO METHODS

Dynamic programming is useful when there is complete knowledge of the environment in the form of transition probabilities. However, most of the time, the environment knowledge is unknown (i.e we known have first hand information about the transition probabilities), and the agent must learn about the environment by interacting with it. In such cases, when knowledge of the environment is missing, we use Monte Carlo methods.

The idea behind Monte Carlo Method in evaluation of the value of states is quite simple. The idea is that, given any state (assuming the robot or agent is in that state), we generate as many trajectories as possible while following the current policy. For instance, in our $3 \times 3$ grid, if the agent is at state $s_3$ and the terminal state is $s_7$, we have illustrated three possible trajectories (red, blue and green). The rewards is $\mathcal{R} = (-1, 3)$

Once we have these trajectories, the next task is to determine the value of each state. With these trajectories, we can only calculate the values of the states that the trajectories pass through. In our $3 \times 3$ example (see Fig. 5), we observe that two trajectories (red and blue) pass through state $s_5$ and all the three trajectories pass through state $s_6$.

To evaluate the value of state $s_5$, we find the return (sum of future rewards) obtained from each trajectory (red and blue) from $s_5$ and to the the terminal state $s_7$.

$$v_{\pi,red}(s_5) = (-1) + (3) = 2$$
$$v_{\pi,blue}(s_5) = (-1) + (3) = 2$$

Finally, we compute the expected value of $s_5$ by averaging the total rewards across all the trajectories that pass through it:

$$v_{\pi,expected}(s_5) = \frac{v_{\pi,red}(s_5) + v_{\pi,blue}(s_5)}{2} = 1$$

Similarly, for the state $s_6$, we have the following

$$v_{\pi,\text{red}}(s_6) = (-1) + (-1) + (3) = 1,$$
$$v_{\pi,\text{blue}}(s_6) = (-1) + (-1) + (3) = 1,$$
$$v_{\pi,\text{green}}(s_6) = (-1) + (-1) + (3) = 1,$$
$$v_{\pi,\text{expected}}(s_6) = \frac{v_{\pi,\text{red}}(s_6) + v_{\pi,\text{blue}}(s_6) + v_{\pi,\text{green}}(s_6)}{3}$$
$$= 1 \tag{29}$$

However, it is important to note that we currently have only three trajectories, each corresponding to a unique episode starting from state $s_3$. In practice, it is essential to conduct numerous episodes originating from all states and identify all trajectories that pass through the specific state whose value we wish to calculate. The more episodes we gather, the greater the likelihood that our estimated value will be close to the true value.

The First-visit and Every-visit Monte Carlo methods in MDP are outlined in Appedix F.

## XI. TEMPORAL DIFFERENCE METHODS

Temporal Difference methods was introduced in 1988 by Sutton [15, 19]. In Temporal Difference methods, the learning process takes place through bootstrapping: we update the value of a state based on the estimated value of the next state. We will explain the idea in details later, but before this let us revisit some equations and derive a relation between them.

Before we move forward, we introduce the following mathematical relation which will be useful in our future derivation. From Eq 15, we have

$$G_t \doteq r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}, \tag{30}$$

It is easy to see the following relation

$$\begin{aligned} G_t &\doteq r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \cdots \\ &= r_{t+1} + \gamma G_{t+1} \end{aligned} \tag{31}$$

Let us revisit Eq. (14) to derive the Bellman equation (23) from it,

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi \left[ G_t \mid S_t = s \right] \quad \text{(From Eq. (14))} \\ &= \mathbb{E}_\pi \left[ r_{t+1} + \gamma G_{t+1} \mid S_t = s \right] \quad \text{(From Eq. (31))} \\ &= \sum_a \pi(a \mid s) \sum_{s' \in \mathcal{S} \setminus \{s\}} \sum_{r \in \mathcal{R}} p(s', r \mid s, a) \\ &\quad \times \left[ r + \gamma \mathbb{E}_\pi \left[ G_{t+1} \mid S_{t+1} = s' \right] \right] \\ &= \sum_a \pi(a \mid s) \sum_{\substack{s' \in \mathcal{S} \setminus \{s_3\} \\ r \in \mathcal{R}}} p(s', r \mid s, a) \\ &\quad \times \left[ r + \gamma v_\pi(s') \right], \quad \text{for all } s \in \mathcal{S}. \tag{32} \end{aligned}$$

where we have used $v_\pi(s') = \mathbb{E}_\pi \left[ G_{t+1} \mid S_{t+1} = s' \right]$

The Bellman equation states that, given a policy $\pi(a \mid s)$ and the environment dynamics $p(s', r \mid s, a)$, we can determine the value of any state $s \in \mathcal{S}$ if the values of all possible subsequent states $s' \in \mathcal{S} \setminus \{s\}$ are known. In other words, if we know the values $v(s')$ for all reachable states $s'$, we can compute the value $v(s)$ by solving the right-hand side of the Bellman equation. For example, in Eq. (27), if $v_\pi(s_2)$ is unknown but the values of other states are known, then using these known state values, we can find $v_\pi(s_2)$ through the Bellman equation. In Eq. (32), an intermediate step is as follows:

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi \left[ G_t \mid S_t = s \right] \\ &= \mathbb{E}_\pi \left[ r_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s \right] \end{aligned} \tag{33}$$

We must note that the action $a$ is determined by the policy structure $\pi$, which defines the agent's behavior in choosing actions. However, the reward $r$ received and the next state $s'$ reached after taking action $a$ are governed by the environment's dynamics. These dynamics are encapsulated in the environment's transition function

$p(s', r \mid s, a)$, which specifies the probability distribution over possible next states and rewards given the current state-action pair.

To interpret Eq. (33), let us consider a deterministic environment.This implies that the agent's trajectory is unique; regardless of the number of episodes, the trajectory remains consistent, eliminating any stochasticity in state transitions. Consequently, we can omit the expectation term from the equation. Under these conditions, the equation suggests that, given the agent follows the policy $\pi$, the value of the current state $s$ is equal to the sum of the reward $r_{t+1}$ obtained after taking action $a$ and the discounted value of the subsequent state $S_{t+1}$ where the agent transitions.

$$v_\pi(s) = r_{t+1} + \gamma v_\pi(S_{t+1}), \quad \text{given} \quad S_t = s \qquad (34)$$

This provides a recursive relationship in which the value of each state can be computed based on the immediate reward and the value of the succeeding state.

Let us change track and introduce an interesting update rule for the state value as follows:

$$v_\pi(s) \leftarrow v_\pi(s) + \alpha \left(r_{t+1} + \gamma v_\pi(S_{t+1}) - v_\pi(s)\right),$$
$$\text{given} \quad S_t = s, \, S_{t+1} = s' \qquad (35)$$

where $\alpha$ is the learning rate.
The update rule in Eq. (35) suggests that the value of a given state $s$ can be adjusted immediately by considering the expected return from the next state $S_{t+1}$ in terms of the immediate reward $r_{t+1}$ and the value $v_\pi(S_{t+1})$. This incremental approach allows for more efficient updates, as it does not require waiting until the end of an episode to evaluate the entire trajectory, in contrast to Monte Carlo methods where updates are only applied at the episode's conclusion. Instead, this rule leverages a technique called temporal difference learning. Temporal difference methods update the value of a state by comparing the current estimate $v_\pi(s)$ with a new estimate based on the observed outcome, specifically $r_{t+1} + \gamma v_\pi(S_{t+1})$. This difference $(r_{t+1} + \gamma v_\pi(S_{t+1}) - v_\pi(s))$ is known as the temporal difference error.

Let us illustrate this using the 2D grid example in Fig. 5. Let the initial value of all the states be zero, $v_\pi(s_1) = v_\pi(s_2) = \ldots = v_\pi(s_9) = 0$. Given policy $\pi$, the agent in state $s = s_3$ moves down to state $s' = s_6$, following the green trajectory. In this process, it receives a reward of $-1$. Thus, the updated value of state $s_3$, following Eq. (35), is given by ($\alpha = 0.1, \gamma = 0.9$)

$$\begin{aligned} v_\pi(s_3) &\leftarrow v_\pi(s_3) + \alpha \left(r_{t+1} + \gamma v_\pi(s_6) - v_\pi(s_3)\right) \\ &\implies v_\pi(s_3) \leftarrow 0 + 0.1 \left(-1 + 0.9 \cdot 0 - 0\right) \\ &\implies v_\pi(s_3) \leftarrow 0 + 0.1 \cdot (-1) \qquad (36) \\ &\implies v_\pi(s_3) \leftarrow -0.1 \\ &\implies v_\pi(s_3) = -0.1 \end{aligned}$$

Let us recall that the primary objective in Temporal Difference learning is to estimate the correct value of each state. Unlike Monte Carlo methods, where the value of a state is updated only at the end of an episode, after observing the full trajectory and accumulating all rewards, Temporal Difference methods update the state value immediately following each action. In here, the update of a state's value happens each time the state is visited, even if revisited within a single trajectory.

For example, in Monte Carlo, the value of a state, say $s_3$, is adjusted only after the episode concludes (follow the green trajectory) and the agent's path has yielded all relevant rewards. In Temporal Difference, however, the update occurs right after each transition, based on the immediate reward and the estimated value of the subsequent state. This approach enables incremental updates across multiple trajectories, leveraging each within-episode transition for faster convergence and improved sample efficiency.

Further, building on the concept of one-step Temporal Difference learning, we can explore two key Temporal Difference based algorithms: SARSA State – Action – Reward – State – Action) and Q-learning. The fundamental difference lies in the update method. SARSA, an on-policy method, updates the value of a state-action pair by considering the actual next action taken by the agent, using the following rule on action value function:

$$q(s, a) \leftarrow q(s, a) + \alpha(r_{t+1} + \gamma q(s', a') - q(s, a)) \qquad (37)$$

Q-learning, an off-policy method, instead updates the value by looking at the highest possible action in the next state, with the rule:

$$q(s, a) \leftarrow q(s, a) + \alpha(r_{t+1} + \gamma \max_{a'} q(s', a') - q(s, a)) \qquad (38)$$

This difference allows SARSA to learn values based on the agent's actual actions, while Q-learning aims for the best possible actions regardless of the policy followed.

## XII. POLICY BASED METHODS

In this section, we will try something new, but our goal stays the same, finding the best possible policy. Let us quickly recall what we did before. Previously, to get an optimal policy, we first found the value of each action. Then, using these action values, we decided which action was best in every state and put together our final policy. But now let us think differently. Can we skip the step of calculating action-values altogether? Surprisingly, the answer is yes.

In this new approach, we follow a different path to discover the optimal policy. Instead of using a fixed policy $\pi(a \mid s)$, where the agent selects actions based on predefined rules or values, we now define the policy as a parameterized function $\pi(a \mid s; \theta)$. The parameter $\theta$ governs the behavior of the policy and can be adjusted through learning. This makes the policy not only flexible but also trainable. By tuning $\theta$, the agent can gradually improve its performance according to experience. Further,

we can extend this formulation to incorporate multiple trainable parameters by introducing a parameter vector $\boldsymbol{\theta} = [\theta_1, \theta_2, \ldots, \theta_n]$, which allows for a more expressive representation of the policy in gradient-based reinforcement learning approaches. In such case, the policy can be represented as $\pi(a \mid s; \boldsymbol{\theta})$.

This method stands in contrast to the value-based approach, where the policy is improved by first estimating action-values and then selecting actions that maximize those values. In the policy gradient framework, the policy is improved directly through changes in the parameters $\boldsymbol{\theta}$, without relying on explicit value estimates. This allows for a more direct and often more scalable way to optimize the agent's behavior, especially in complex or continuous action spaces.

### A. Policy Gradient I: The Theory

In policy gradient methods, our goal is to directly optimize the performance of a parameterized policy, $\pi(a \mid s; \boldsymbol{\theta})$. Let us define a scalar objective function $J(\boldsymbol{\theta})$ that measures the quality of the policy. Specifically, $J(\boldsymbol{\theta})$ is defined as the *expected return* when the agent starts interacting with the environment at the initial step $t = 0$ and follows the policy $\pi(a \mid s; \boldsymbol{\theta})$ (in short notation $\pi_{\boldsymbol{\theta}}$) thereafter,

$$J(\boldsymbol{\theta}) \doteq \mathbb{E}_{\pi_{\boldsymbol{\theta}}}[G_t], \quad \text{with } t = 0. \tag{39}$$

remember that we expected return because it provides a stable and reliable measure of a policy's average performance in a stochastic environment (refer back to State Value Function subsection in Sec. IV).

Let us understand how Eq. (39) is different from our earlier definition $v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$ in Eq. (14). The key distinction lies in the conditioning of the expectation. The expression $v_\pi(s)$ represents the expected return starting from a specific state $s$, and thus captures how desirable it is to be in that particular state under policy $\pi$. In contrast, $J(\boldsymbol{\theta})$ evaluates the expected return starting from the initial time step $t = 0$. Therefore, $J(\boldsymbol{\theta})$ provides a global measure of the overall performance of the policy, averaged across entire episodes.

Now that policy $\pi_{\boldsymbol{\theta}}$ is defined as a differentiable function of parameters $\boldsymbol{\theta}$ and the objective $J(\boldsymbol{\theta})$, policy gradient methods aim to improve policy by maximizing $J(\boldsymbol{\theta})$. This is done by computing the gradient $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$, which indicates the direction in which the parameters should be adjusted to increase the expected return. Say, $\boldsymbol{\theta}$ is a vector of $d$ parameters, then the gradient $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ is itself a vector, given by

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix}. \tag{40}$$
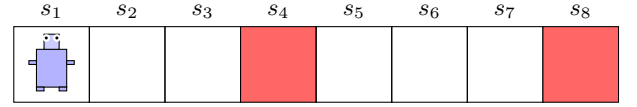


FIG. 6: Grid world with robot in $s_1$ and terminal states $s_4$ and $s_8$ highlighted in red

The parameters are then updated using gradient ascent:

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta} = \boldsymbol{\theta}_k}, \tag{41}$$

where $\alpha$ is the learning rate. In this way, the policy is directly optimized without requiring a value function estimate.

Note that, the gradient $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ is a vector of partial derivatives, one for each parameter $\theta_i$. Thus, the update rule in Eq. (41) adjusts all parameters simultaneously:

$$\theta_{i,k+1} = \theta_{i,k} + \alpha \left. \frac{\partial J}{\partial \theta_i} \right|_{\boldsymbol{\theta} = \boldsymbol{\theta}_k}, \quad \text{for } i = 1, \ldots, d. \tag{42}$$

In other words, every parameter is pulled in the direction that increases the objective, and the step size is controlled by the learning rate $\alpha$.

### B. Policy Gradient II: Example

Let us understand how we can implement the formulation of the policy gradient for the 1D grid example (Fig. 7). The example in Fig. 1 was easy for the agent to learn. The robot just needed to go left if the terminal state was on its left, and right if it was on its right. The new example in Fig. 7 is more challenging because it has two terminal states $s_4$ and $s_8$. In state $s_5$, the robot should always go right. In state $s_6$, both left and right lead to terminal states, so the agent needs to choose wisely. In state $s_7$, the best action is again to go right. This setup shows that learning becomes harder when the environment has more complex decisions.

If we were to implement Monte Carlo or Temporal difference method, we could start with the following random policy:

```
policy = {
    1: {'right': 1.0},
    2: {'left': 0.5, 'right': 0.5},
    3: {'left': 0.5, 'right': 0.5},
    4: {},
    5: {'left': 0.5, 'right': 0.5},
    6: {'left': 0.5, 'right': 0.5},
    7: {'left': 0.5, 'right': 0.5},
    8: {}
}
```

The initial random policy must be improved to reach a stable optimal policy. However, since we are working within the policy gradient framework, it is essential to

define a trainable parameter, denoted by $\theta$, or more generally as a parameter vector $\boldsymbol{\theta} = [\theta_1, \theta_2, \ldots, \theta_6]$. Consequently, we must revise our initial policy representation from $\pi(a \mid s)$ to the parameterized form $\pi(a \mid s; \boldsymbol{\theta})$ to reflect this dependency. A simple way to do this is as follows (one must remember that the dependency on $\theta$ can be chosen in a more complex form):

$$\pi(a \mid s; \theta) = \begin{array}{c|c|c} s & \pi(left \mid s;\theta) & \pi(right \mid s;\theta) \\ \hline 1 & 0 & 1 - \theta \\ 2 & 0.5 + \theta & 0.5 - \theta \\ 3 & 0.5 + \theta & 0.5 - \theta \\ 4 & 0 & 0 \\ 5 & 0.5 + \theta & 0.5 - \theta \\ 6 & 0.5 + \theta & 0.5 - \theta \\ 7 & 0.5 + \theta & 0.5 - \theta \\ 8 & 0 & 0 \end{array}$$

or in more general form

$$\pi(a \mid s; \theta) = \begin{array}{c|c|c} s & \pi(left \mid s;\theta) & \pi(right \mid s;\theta) \\ \hline 1 & 0 & 1 - \theta_1 \\ 2 & 0.5 + \theta_2 & 0.5 - \theta_2 \\ 3 & 0.5 + \theta_3 & 0.5 - \theta_3 \\ 4 & 0 & 0 \\ 5 & 0.5 + \theta_4 & 0.5 - \theta_4 \\ 6 & 0.5 + \theta_5 & 0.5 - \theta_5 \\ 7 & 0.5 + \theta_6 & 0.5 - \theta_6 \\ 8 & 0 & 0 \end{array}$$

The general policy gradient update is given by Eq. (41), where the goal is to adjust the policy parameters $\boldsymbol{\theta}$ in the direction that maximizes the expected return $J(\boldsymbol{\theta})$. However, directly computing $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ is difficult because it involves expectations over all possible trajectories. To address this, the REINFORCE algorithm[20] uses a clever trick: it approximates this gradient using sampled trajectories and the log-likelihood trick. Specifically, it substitutes the gradient with

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \approx G_t \cdot \nabla_{\boldsymbol{\theta}} \log \pi(a_t \mid s_t; \boldsymbol{\theta}), \qquad (43)$$

where $G_t$ is the return from the time step $t$ onward. This leads to the *reinforce* update rule:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \cdot G_t \cdot \nabla_{\boldsymbol{\theta}} \log \pi(a_t \mid s_t; \boldsymbol{\theta}). \qquad (44)$$

This formulation allows learning from actual sampled episodes by increasing the probability of actions that lead to higher returns, even without knowing the full model of the environment.

Let us now look at the improved policy for Fig. 7,

```
improved_policy = {
"s1": {'right': 1.00, 'left': 0.00},
"s2": {'right': 0.99, 'left': 0.01},
"s3": {'right': 0.99, 'left': 0.01},
"s4": {},  # Terminal state
"s5": {'right': 0.01, 'left': 0.99},
```
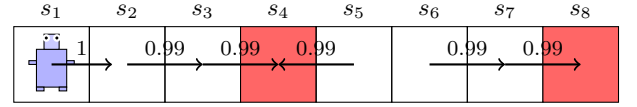


FIG. 7: Improved policy with action directions and probabilities derived from the policy gradient method. Terminal states $s_4$ and $s_8$ are highlighted in red. The numerical values indicate the transition probabilities for each action, demonstrating the learned stochastic policy after convergence.

```
"s6": {'right': 0.99, 'left': 0.01},
"s7": {'right': 0.99, 'left': 0.01},
"s8": {}  # Terminal state
}
```

## XIII.   ACTOR-CRITIC METHOD

Policy gradient methods, as introduced earlier in Sec. XI, works quite well in environments where rewards are sparse (agent only receive reward at the end of an episode), such as the 1D grid world. In those settings, the agent learns by seeing the total result of an episode, which helps guide its future decisions.

However, for long episodes in sparse reward settings, one of the critical issues is that they suffer from high variance in their gradient estimates. A simple policy gradient algorithm updates its policy based on the total return of an entire episode. The problem is that the episode's total reward can vary wildly due to randomness, even if the agent's actions are similar. This makes the gradient updates noisy and unstable, leading to slow and inefficient learning.[21–23].

Also, policy gradient is a very slow learning method like Monte Carlo.[24, 25]. This is because they typically wait until the entire episode is complete before making any updates. And that is why, in environments with long tasks or many steps (known as long-horizon environments), this leads to inefficient and unstable training.

To solve these problems, the *actor-critic method*[26] offers a useful solution. It combines the strengths of both value-based methods (like Temporal Difference) and policy-based methods (like policy gradient). As the name suggests:

- The *actor* is the part that *decides what to do* it updates the policy $\pi(a|s;\theta)$, just like in regular policy gradient.

- The *critic* is the part that *evaluates how good a situation is*—it uses a value function $V(s;w)$, learned using temporal difference techniques.

The policy gradient update rule (introduced in Sec. XI) is modified as follows:

$$\nabla_{\theta} J(\theta) \approx A(s_t, a_t) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t) \qquad (45)$$

where:

- $J(\theta)$ is the performance objective, i.e., the expected total reward.

- $\pi_\theta(a_t \mid s_t)$ is the probability of taking action $a_t$ in state $s_t$ under the policy parameterized by $\theta$.

- $\nabla_\theta \log \pi_\theta(a_t \mid s_t)$ represents how a small change in $\theta$ affects the log-probability of the chosen action.

- $A(s_t, a_t)$ is the *advantage function*, which measures how much better or worse action $a_t$ was compared to the critic's baseline expectation for state $s_t$.

In the actor–critic algorithm, the advantage function $A(s_t, a_t)$ serves as a replacement for the return $G_t$ used in Equation. (43) of the REINFORCE policy gradient method. The advantage function is often estimated using the *temporal difference error*:

$$A(s_t, a_t) \approx \delta_t = r_{t+1} + \gamma V(s_{t+1}; w) - V(s_t; w), \quad (46)$$

where:

- $r_{t+1}$ is the immediate reward received after taking action $a_t$ in state $s_t$,

- $\gamma \in [0, 1]$ is the discount factor determining the importance of future rewards,

- $V(s_t; w)$ is the critic's estimate of the value of state $s_t$,

- $V(s_{t+1}; w)$ is the critic's estimate of the value of the next state.

The temporal difference error shows how surprising the outcome is. If the reward we receive, along with the value of the next state, is greater than expected, then $\delta_t > 0$. If it is smaller than expected, then $\delta_t < 0$.

So now there will be two updates. The critic updates its value function parameters $w$ according to:

$$w \leftarrow w + \beta\, \delta_t\, \nabla_w V(s_t; w), \qquad (47)$$

where $\beta > 0$ is the critic's learning rate, and $\nabla_w V(s_t; w)$ is the gradient of the value function with respect to its parameters. This update moves the critic's predictions closer to the observed outcomes.

Finally, the actor updates its policy parameters $\theta$ using the same temporal difference error:

$$\theta \leftarrow \theta + \alpha\, \delta_t\, \nabla_\theta \log \pi_\theta(a_t \mid s_t), \qquad (48)$$

where $\alpha > 0$ is the actor's learning rate.

**Actor-Critic Update Example on 1D Grid**

Let us take the 1D grid example. Initially we have,

```
π(a | s;θ) = {
        1: {right:  1 − θ₁},
        2: {left:  0.5 + θ₂, right:  0.5 − θ₂},
        3: {left:  0.5 + θ₃, right:  0.5 − θ₃},
        4: {},
        5: {left:  0.5 + θ₄, right:  0.5 − θ₄},
        6: {left:  0.5 + θ₅, right:  0.5 − θ₅},
        7: {left:  0.5 + θ₆, right:  0.5 − θ₆},
        8: {}
         }
```

All policy parameters are initialized to zero: $\theta_1 = \theta_2 = \cdots = \theta_6 = 0$

Here learning rate $\alpha = 0.1$ (for both actor and critic update) and discount factor $\gamma = 0.9$. The initial value function is $v(s) = 0$ for all states.

Consider the trajectory $s_6 \to s_7 \to s_8$ (terminal). The temporal difference error and critic update are:

$$\delta_t = r_t + \gamma v(s_{t+1}) - v(s_t), \quad v(s_t) \leftarrow v(s_t) + \alpha \delta_t \quad (49)$$

Transition $s_6 \to s_7$ (reward $-1$):

$$v(s_6) = 0, \quad v(s_7) = 0$$
$$\delta = -1 + 0.9 \cdot 0 - 0 = -1$$
$$v(s_6) \leftarrow 0 + 0.1 \cdot (-1) = -0.1$$

Transition $s_7 \to s_8$ (reward $+5$, terminal $s_8$):

$$v(s_7) = 0, \quad v(s_8) = 0$$
$$\delta = 5 + 0.9 \cdot 0 - 0 = 5$$
$$v(s_7) \leftarrow 0 + 0.1 \cdot 5 = 0.5$$

After critic updates:

$$v(s_6) = -0.1, \quad v(s_7) = 0.5, \quad v(s_8) = 0$$

Next step is to calculate the advantage function,

$$A(s_t, a_t) = r_t + \gamma v(s_{t+1}) - v(s_t) \qquad (50)$$

- For $s_6 \to s_7$:

$$A(s_6, right) = -1 + 0.9 \cdot 0.5 - (-0.1) = -0.45$$

- For $s_7 \to s_8$:

$$A(s_7, right) = 5 + 0.9 \cdot 0 - 0.5 = 4.5$$

The actor update is:

$$\theta \leftarrow \theta + \alpha\, A(s_t, a_t)\, \nabla_\theta \log \pi_\theta(a_t \mid s_t) \qquad (51)$$

With the linear parametrization:

$$\pi(right \mid s) = 0.5 - \theta, \qquad \pi(left \mid s) = 0.5 + \theta$$

the derivative is:

$$\frac{\partial}{\partial \theta} \log \pi(right \mid s) = -\frac{1}{0.5 - \theta}$$

Update at $s_6$ (action = right, $\theta_6 = 0$):

$$\Delta\theta_6 = \alpha\, A(s_6, right) \left(-\frac{1}{0.5 - \theta_6}\right)$$
$$= 0.1 \times (-0.45) \times (-2)$$
$$= 0.09$$
$$\theta_6' = 0 + 0.09 = 0.09$$
$$\pi'(right \mid s_6) = 0.5 - 0.09 = 0.41$$
$$\pi'(left \mid s_6) = 0.5 + 0.09 = 0.59$$

Update at $s_7$ (action = right, $\theta_7 = 0$):

$$\Delta\theta_7 = 0.1 \times 4.5 \times (-2) = -0.9$$
$$\theta_7' = 0 - 0.9 = -0.5 \quad \text{(clipped to } [-0.5, 0.5])$$
$$\pi'(right \mid s_7) = 0.5 - (-0.5) = 1.0$$
$$\pi'(left \mid s_7) = 0.0$$

Final Values and Policy:

$$v(s_6) = -0.1, \quad v(s_7) = 0.5, \quad v(s_8) = 0$$
$$\pi(right \mid s_6) = 0.41, \quad \pi(left \mid s_6) = 0.59$$
$$\pi(right \mid s_7) = 1.0, \quad \pi(left \mid s_7) = 0.0$$

These updates illustrate how the actor-critic framework adjusts the value estimates and policy probabilities step by step along a trajectory in a 1D grid world.

## XIV. QUANTUM CONTROL: INTRODUCTION

The concept of quantum control emerged in the mid-1980s, when scientists realized that it might be possible to deliberately guide the behavior of quantum systems using carefully designed electromagnetic fields. In quantum mechanics, atoms and molecules can reach the same final state through multiple quantum pathways, and these pathways can interfere with each other much like overlapping ripples on water. Researchers such as Brumer and Shapiro (1986) demonstrated that by shaping the phase and timing of laser pulses, one could exploit this interference to control chemical reactions and molecular transitions [27]. This insight laid the foundation for what became known as *coherent quantum control*. Subsequent theoretical and experimental studies expanded this idea into a general framework for manipulating atomic and molecular dynamics [28, 29].

At first, the theoretical framework remained largely conceptual because probing and manipulating quantum systems at such fine scales was extremely difficult. These ideas could not be realized experimentally until the 1990s, which marked a period of rapid experimental progress enabled by the development of *femtosecond laser technology* and *pulse-shaping techniques*. These advances made it possible to synthesize and control ultrashort laser pulses on timescales comparable to molecular vibrations and electronic transitions, transforming coherent control from a theoretical concept into an experimentally realizable discipline [30, 31].

At the same time, significant progress was made on the theoretical and computational side of quantum control. Researchers developed *Quantum Optimal Control Theory (QOCT)*, a framework that combines the time-dependent Schrödinger equation with optimization methods to design external control fields that achieve specific goals, such as transferring population between states, implementing high-fidelity quantum gates, or reducing decoherence [32]. A major milestone was the introduction of the GRAPE algorithm (Gradient Ascent Pulse Engineering) by Khaneja et al. (2005), which provided a fast and efficient way to solve high-dimensional optimization problems. Originally developed for nuclear magnetic resonance (NMR) systems, GRAPE was later adapted for use in quantum information processing [33]. Since then, improved variants such as Krotov's method[34], CRAB (Chopped Random Basis)[35], and GOAT (Gradient Optimization of Analytic Controls)[36] have further extended the computational tools available for designing precise quantum control fields [37].

By the 2010s, research in quantum control increasingly focused on quantum technologies, including superconducting qubits, trapped ions, and solid-state spin systems. In these platforms, highly accurate and noise-resilient control is essential for building high-fidelity quantum gates and achieving scalable quantum computation [38]. The growing complexity of large, open quantum systems exposed to environmental noise motivated researchers to explore machine learning (ML) approaches to automate and improve control design. ML-based quantum control leverages data-driven models, neural networks, and gradient-based optimization to learn control strategies that are adaptive, robust, and efficient for practical hardware [39].

The field of Quantum Reinforcement Learning(QRL) combines ideas from control theory, quantum dynamics, and artificial intelligence. In QRL, controlling a quantum system is formulated as a sequential decision-making process. We have already learned the components of RL, in this case the quantum system acts as the environment, the applied control pulses represent the actions, and performance metrics such as fidelity or energy cost serve as the reward function. Reinforcement learning (RL) agents can therefore learn optimal control strategies either by interacting directly with experiments (model-free) or by training on simulations of quantum dynamics (model-based).

Recent studies have demonstrated that deep RL can effectively perform tasks such as quantum state preparation, gate optimization, and error suppression even in the presence of environmental noise and decoherence [40–

46]. With continued progress in quantum hardware and hybrid quantum-classical computing, QRL is a promising direction toward autonomous, adaptive, and robust controllers for next-generation quantum technologies.

## XV. APPLICATION OF RL IN QUANTUM CONTROL

The goal of optimal terminal control is to drive the system from a known initial quantum state to a desired target state at a specified terminal (i.e., final) time $T$. The evolution of the quantum system is governed by the time-dependent Schrödinger equation:

$$i\hbar\frac{d\,|\psi(t)\rangle}{dt} = H(u(t))\,|\psi(t)\rangle,\qquad(52)$$

where $|\psi(t)\rangle$ is the quantum state at time $t$, and $H(u(t))$ is the Hamiltonian, which depends on the control field $u(t)$.

To quantify the success of the control, we use the *fidelity* between the final state $|\psi(T)\rangle$ and the desired target state $|\psi_{\text{target}}\rangle$:

$$F = |\langle\psi_{\text{target}}|\psi(T)\rangle|^2.\qquad(53)$$

A fidelity $F = 1$ indicates perfect overlap with the target state.

### A. Example problem in Optimal control

Consider a single qubit, initially in state $|0\rangle$, and we aim to drive it to the target state $|1\rangle$ using a control field. The control is applied through a rotation operator along the X-axis, which is parameterized by a rotation angle $\theta$.

- Initial state: $|0\rangle$

- Target state: $|1\rangle$

- Control parameter: $\theta$

The unitary operator for the X-axis rotation is given by:

$$U(\theta) = e^{-i\frac{\theta}{2}\sigma_x}\qquad(54)$$

where $\sigma_x$ is the Pauli-X operator.

#### Using Actor-critic Method

We solve this example by actor-critic method discussed earlier. The actor decides what action to take based on the current state, while the critic evaluates how good the action taken by the actor is by estimating the value function. Let's see how to integrate it all in RL:

- **State** $s$: The quantum state of the qubit, represented as a 2D vector in the Bloch sphere(See fig 8). For example, $|0\rangle$ is $[1, 0]$ and $|1\rangle$ is $[0, 1]$.

- **Action** $a$: The rotation angle $\theta$ applied to the qubit.

- **Reward** $R(s, a)$: The fidelity between the final state and the target state, defined as:

$$R(s, a) = |\langle\psi_{\text{target}}|\psi_{\text{final}}\rangle|^2\qquad(55)$$

- **Policy**: The agent's strategy, represented by the actor network, which outputs the optimal rotation angle $\theta$.

- **Value function**: The critic network estimates the expected future reward for each state.

In our quantum control problem, the actor network learns to output the optimal rotation angle $\theta$, while the critic network estimates the expected cumulative reward (derived from fidelity) for each state. The actor and critic are trained using the following rules: :

- The actor's objective is to maximize the expected reward by adjusting $\theta$. Its loss function is defined as

$$L_{\text{actor}} = -A(s, a) \cdot \log \pi_\theta(a|s),\qquad(56)$$

where $\pi_\theta(a|s)$ is the probability of taking action $a$ in state $s$ under the actor policy, and $A(s, a)$ is the advantage, i.e., the difference between the actual reward and the critic's value estimate.

- The critic is trained to minimize the temporal-difference error:

$$L_{\text{critic}} = \left(R_{\text{actual}} - V_{\text{estimated}}\right)^2,\qquad(57)$$

where $R_{\text{actual}}$ is the reward observed from the environment, and $V_{\text{estimated}}$ is the value function predicted by the critic.

After training the actor and critic networks, we evaluated the learned control strategy by computing the fidelity between the final quantum state and the target state. The fidelity is given by the overlap between the final state $|\psi_{\text{final}}\rangle$ and the target state $|\psi_{\text{target}}\rangle$:

$$F = |\langle\psi_{\text{target}}|\psi_{\text{final}}\rangle|^2\qquad(58)$$

The trained agent successfully learned the optimal rotation angle to drive the qubit from the initial state $|0\rangle$ to the target state $|1\rangle$ as shown in the figure 8. The final fidelity achieved was $F = 0.99999998596$, indicating a high degree of accuracy in reaching the target state.

This problem serves primarily as a pedagogical example, and while it effectively illustrates key concepts, employing an actor-critic algorithm in this context may be considered disproportionate given the simplicity of the task.
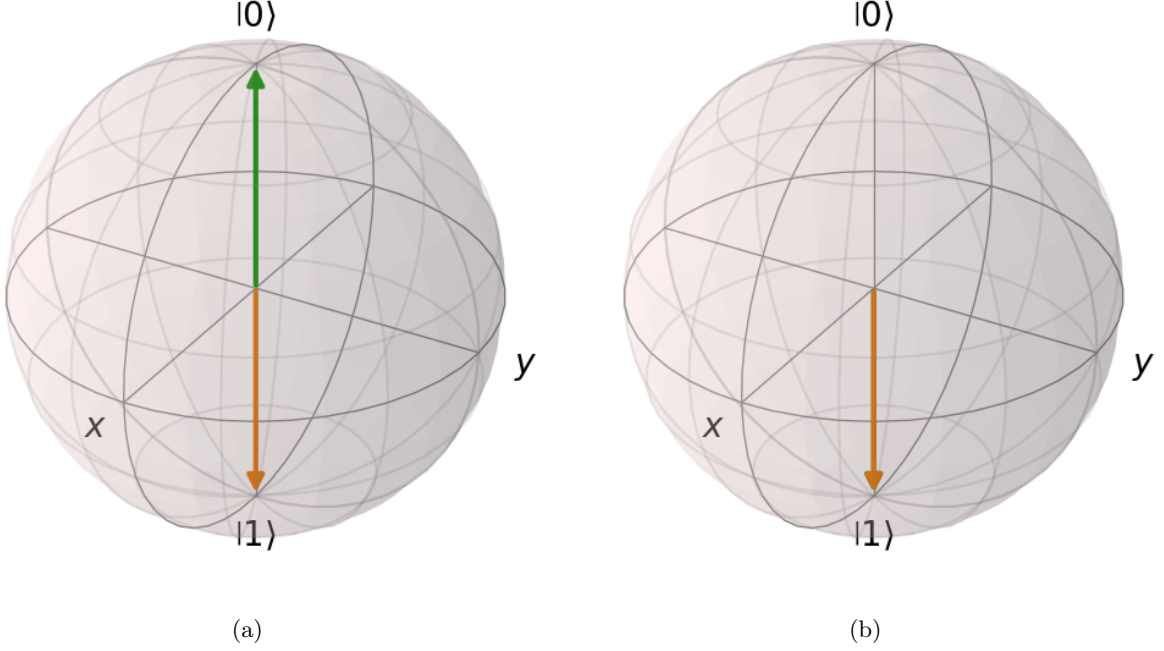
FIG. 8: Bloch sphere representations of quantum state evolution. (**a**) The initial state (in green) corresponds to the computational basis state $|0\rangle$, and the target state (in orange) is $|1\rangle$. (**b**) The final state after control (in green) overlaps with the target state (in orange), indicating successful evolution. The overlap causes only one visible vector due to their alignment.

## B. Future Direction: Tracking control in Quantum systems

Till now, we have discussed RL primarily in the context of quantum optimal control. A growing body of work has demonstrated that RL is highly effective for optimal control tasks such as quantum state preparation [47], gate synthesis [48], and dynamical stabilization [49], achieving near-unit fidelities even in complex and noisy settings.

In most of these studies, the control objective is formulated in terms of a final time target, and the learning agent is rewarded based on a terminal fidelity or cost. As a result, the intermediate system dynamics are optimized only indirectly, insofar as they help reach the desired final state.

There are, however, many physically relevant situations where this formulation is not sufficient. In such cases, the control task requires not only reaching the correct final state, but also regulating the system's behavior throughout its evolution. This naturally leads to the concept of tracking control, where the objective is to guide the system along a prescribed time dependent trajectory rather than optimizing a terminal outcome.

Early developments in quantum tracking control, particularly in the early 2000s, were based on analytical constructions [50, 51]. These approaches were inspired by inverse dynamics and exact tracking methods developed earlier in classical nonlinear control theory [52, 53]. The central idea was to derive the control field explicitly from the equations of motion such that a selected observable follows a predefined trajectory.

Concretely, consider a controlled quantum system governed by the Hamiltonian

$$H(t) = H_0 + u(t)H_c, \tag{59}$$

where $u(t)$ denotes the control field. If $O$ is the observable of interest, the tracking objective is to enforce

$$y(t) = \langle O(t) \rangle \approx y_d(t), \tag{60}$$

where $y_d(t)$ is a prescribed target trajectory. The control task is therefore defined over the full time interval of evolution, rather than at a single final time.

The time derivative of the observable expectation value can be written as

$$\frac{d}{dt}\langle O \rangle = \frac{i}{\hbar}\langle [H_0, O] \rangle + \frac{i}{\hbar}u(t)\langle [H_c, O] \rangle. \tag{61}$$

This expression highlights a key structural feature. The observable dynamics depend linearly on the control field. If a desired trajectory $y_d(t)$ is specified, one can impose

$$\frac{d}{dt}\langle O \rangle = \dot{y}_d(t) \tag{62}$$

and formally solve for the control field $u(t)$ by inverting the above relation. In this way, analytical tracking enforces the desired trajectory by construction.

While such analytical tracking schemes are elegant and physically transparent, they also reveal important practical limitations. The inversion can become singular, it

relies on accurate knowledge of the system Hamiltonian, and it becomes increasingly fragile in the presence of decoherence, noise, or high-dimensional dynamics. Consequently, analytical tracking is most effective for few-level, well characterized systems.

Tracking control can be understood as an inverse-control problem, namely determining the control action required to produce a desired dynamical response. This perspective provides the natural entry point for RL.

In RL-based tracking control, the inverse mapping no longer needs to be written down explicitly. Instead, the control field is generated by a policy,

$$u(t) = \pi(s(t)),\tag{63}$$

where the policy $\pi$ maps the agent's state $s(t)$ to a control action. For tracking problems, the state typically encodes the tracking error,

$$e(t) = y(t) - y_d(t),\tag{64}$$

and possibly its recent history or local temporal variation. Rather than enforcing $e(t) = 0$ instantaneously, the agent can be trained to reduce the error over time.

This can be achieved by defining a reward signal that penalizes deviations from the desired trajectory, for example,

$$r(t) = -e^2(t).\tag{65}$$

The learning objective then becomes the minimization of the accumulated tracking error over the full evolution,

$$\int_0^T e^2(t)\, dt,\tag{66}$$

rather than the optimization of a terminal fidelity. This shift from instantaneous inversion to trajectory level optimization is a defining feature of RL based tracking control.

Within this setting, several classes of RL algorithms are proposed below, each addressing different aspects of the tracking problem.

A natural starting point is policy-gradient methods such as [54], Trust Region Policy Optimization (TRPO) [55], and Proximal Policy Optimization (PPO) [56]. In these approaches, the control field is generated directly by a parameterized policy,

$$u(t) = \pi_\theta(s(t)),\tag{67}$$

and the policy parameters are updated by maximizing the expected cumulative reward. For tracking problems, the reward is typically defined in terms of the instantaneous tracking error (Eq.65). Because policy gradient methods optimize trajectory level objectives without requiring an explicit system model, they are well suited for enforcing time-dependent behavior and smooth control fields. PPO is particularly attractive due to its numerical stability and ease of implementation. Another naturally aligned class of methods is actor-critic reinforcement learning. Algorithms such as DDPG [57], TD3 [58], and SAC [59] may be promising candidates for future tracking control applications. DDPG may be suitable for deterministic tracking control, TD3 could offer improved stability for long-time tracking, and SAC may enable robust tracking by incorporating entropy, which can be beneficial in the presence of noise and uncertainty.

RL therefore, offers a natural framework for future developments in tracking control of quantum systems, enabling model-free optimization of time-dependent objectives with stable and experimentally feasible control policies.

## XVI. QUANTUM ADVANTAGE IN LEARNING: QUANTUM REINFORCEMENT LEARNING

### A. Quantum Computing: Introduction

Any computing is ultimately done using some rules of nature. We first discover the rules/laws that govern how the physical world behaves, and then we build devices that obey those laws to help us compute. Every computing system, whether mechanical, electrical, or quantum, is an embodiment of the rules of nature translated into hardware and logic.

In the beginning, we learned to harness electricity. From its behavior, engineers formulated a set of simple and reliable rules, and from those rules we built all of the classical computing. The entire digital world rests on the deterministic physics of electric charge and current. The following set of principles constitute the axiomatic framework underlying classical digital computation.

1. *Binary State Representation*: Electrical signals exist in one of two discrete voltage levels, encoding logical states $\{0, 1\}$.

2. *Deterministic Computation*: Circuit outputs are uniquely determined by their inputs, exhibiting no inherent stochasticity.

3. *Boolean Logic Operations*: Computational primitives (AND, OR, NOT) obey the laws of Boolean algebra.

4. *Non-destructive Measurement*: Information can be read and replicated without altering the original state.

5. *Physical Parallelism*: Concurrent operations require proportional physical resources (additional circuits or processors).

Every bit of data that moves through a computer today, every image, word, and number, is a pattern of voltages traveling through circuits that obey these classical electrical laws. However, when researchers looked deeper into nature, beneath wires and transistors, they found

that these rules were only approximations of a more profound layer of reality, the quantum reality of nature. Thus came a new set of rules, the rules of quantum mechanics, which is now the basis of quantum computation. Where classical computers manipulate *bits* constrained by electrical laws, quantum computers process *qubits*, quantum systems that obey the quantum rules.

Let us now summarize the quantum rules:

1. *Quantum Superposition*: A qubit exists as a linear combination of basis states: $|\psi\rangle = \alpha|0\rangle + \beta|1\rangle$, where $\alpha, \beta \in \mathbb{C}$.

2. *Measurement-Induced Collapse*: Observation projects the qubit to $|0\rangle$ or $|1\rangle$ with probabilities $|\alpha|^2$ and $|\beta|^2$, respectively, irreversibly altering the quantum state.

3. *Unitary Evolution*: Quantum operations are reversible transformations in a Hilbert space, preserving total information.

4. *Quantum Entanglement*: Multiple qubits can form inseparable joint states, exhibiting non-classical correlations.

5. *No-Cloning Theorem*: Unknown quantum states cannot be duplicated; information is transformed but not copied.

6. *Quantum Interference*: Computational paths are explored in superposition, with constructive and destructive interference amplifying correct solutions.

The main advantage of quantum computing lies in its ability to outperform classical systems on certain problems, achieving exponential speedups. In simple terms, superposition allows a quantum computer to explore many possible inputs at once, rather than one by one as classical computers do.

### B.   Quantum Computing: Technical Details

Let us now get into technical details of how quantum computing works. A quantum circuit is the basic blueprint of how computation happens in a quantum computer. A classical circuit uses wires and logic gates to manipulate bits (0s and 1s) whereas a quantum circuit uses quantum gates to manipulate qubits. Let us begin with a simple qubit, which can exist in a combination of two basic states, $|0\rangle$ and $|1\rangle$:

$$|\psi\rangle = \alpha|0\rangle + \beta|1\rangle, \qquad (68)$$

Quantum gates implement unitary transformations that evolve the quantum state:

$$|\psi'\rangle = U|\psi\rangle, \quad U^\dagger U = UU^\dagger = I. \qquad (69)$$

Single-qubit gates perform rotations on the Bloch sphere, parameterized as:

$$R_j(\theta) = \exp\left(-i\frac{\theta}{2}\sigma_j\right), \quad j \in \{x, y, z\}, \qquad (70)$$

where $\sigma_j$ denotes the Pauli matrices. Multi-qubit gates, such as the controlled-NOT (CNOT), create entanglement between qubits, a uniquely quantum resource absent in classical computation.

The measurement process, repeated over many trials (also called shots), yields probability distributions or expectation values:

$$\langle O \rangle = \langle \psi | U^\dagger O U | \psi \rangle, \qquad (71)$$

where $O$ represents an observable operator (e.g., Pauli-$Z$).

In practice, a single measurement of a quantum circuit yields a random outcome, such as `0110`. The power of quantum measurement lies not in one result but in the pattern revealed through many repetitions. Each run, or shot, is like a single exposure in photography, unclear alone, but sharp in aggregate. By repeating the circuit, we estimate probabilities or expectation values like

$$\langle O \rangle = \langle \psi | U^\dagger O U | \psi \rangle, \qquad (72)$$

where $O$ is an observable such as the Pauli-$Z$ operator. Thus, quantum measurement uncovers information statistically.

### C.   Quantum Reinforcement Learning

RL began as a behavioral concept rooted in psychology and control theory, which describes how organisms learn through trial and error by responding to rewards and punishments [60]. This intuitive idea gradually evolved into a rigorous computational framework in which an agent interacts with an environment to maximize cumulative rewards [19, 60]. The modern era of machine learning has reshaped RL into a model-driven approach. Agents are now represented by neural networks and other parameterized systems that learn to make better decisions through experience [60]. What once described animal learning has become a core framework for developing intelligent, adaptive machines.

This computational framework now extends into the quantum realm. Quantum Reinforcement Learning (QRL) [61–63] is one of the three principal branches of Quantum Machine Learning (QML), alongside Quantum Supervised and Quantum Unsupervised Learning. Like classical ML, QRL is built around models characterized by trainable parameters—adjustable quantities that enable the system to capture patterns and improve performance through training. A model endowed with such adaptive parameters is known as a learnable model.

In QRL, learnable system is realized through parameterized quantum circuits, where the rotation angles $\theta$ in

gates such as $R_x(\theta)$, $R_y(\theta)$, and $R_z(\theta)$ act as trainable parameters. Such circuits, known as Variational Quantum Circuits (VQCs), form the foundation of learnable quantum models. By iteratively tuning $\theta$ via classical optimization, the circuit adapts its behavior to perform computational or decision-making tasks. This concept lies at the heart of Quantum Machine Learning [64, 65].

To understand quantum RL in more detail, let us revisit policy gradient in Sec. XII. In the quantum setting, the policy $\pi(a \mid s; \theta)$ is realized by a variational quantum circuit (VQC), and the trainable vector $\theta$ consists of the gates' rotation angles. When the policy is a VQC, the terms $\nabla_\theta \log \pi(a_t \mid s_t; \theta)$ in Eq 43 are obtained on the quantum model, while the optimizer runs classically. Concretely, a state is encoded into qubits, the circuit is executed, measurements yield action probabilities, and an action is sampled. After observing returns, those angles are updated using the standard REINFORCE policy-gradient estimator, where gradients are computed via the parameter-shift rule on the quantum circuit and optimization runs classically. In short, the loop is exactly the classical policy-gradient algorithm; the only change is that the function class for $\pi$ is a quantum circuit whose learnable rotation angles form $\theta$.

## XVII. CONCLUSION

In this tutorial, we have explained several key concepts of Reinforcement Learning (RL) through simple and intuitive examples, making it easy for students to grasp the foundational topics. Most of the fundamental concepts in RL are covered with straightforward, easy-to-follow tutorials. To complement the theoretical explanations, we have included code implementations that students can explore. Supporting materials, such as booklets, provide additional guidance, while the complete implementations of the examples discussed in this tutorial are available on GitHub for reference and hands-on practice.

[1] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, Timothy Lillicrap, and David Silver. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, December 2020.

[2] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, December 2018.

[3] Marc G. Bellemare, Salvatore Candido, Pablo Samuel Castro, Jun Gong, Marlos C. Machado, Subhodeep Moitra, Sameera S. Ponda, and Ziyu Wang. Autonomous navigation of stratospheric balloons using reinforcement learning. *Nature*, 588(7836):77–82, December 2020.

[4] Ali Irshayyid, Jun Chen, and Guojiang Xiong. A review on reinforcement learning-based highway autonomous vehicle control. *Green Energy and Intelligent Transportation*, 3(4):100156, August 2024.

[5] Md. Al-Masrur Khan, Md Rashed Jaowad Khan, Abul Tooshil, Niloy Sikder, M. A. Parvez Mahmud, Abbas Z. Kouzani, and Abdullah-Al Nahid. A systematic review on reinforcement learning-based robotics within the last decade. *IEEE Access*, 8:176598–176623, 2020.

[6] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022.

[7] Abhijit Sen Sonali Panda. Reinforcement tutorial for undergraduates — booklets. https://github.com/asen009/Reinforcement_Tutorial_Undergraduates/tree/main/Booklets, 2025. GitHub repository.

[8] Bernhard Jaeger and Andreas Geiger. An invitation to deep reinforcement learning, 2023.

[9] Abhijit Gosavi. Reinforcement learning: A tutorial survey and recent advances. *INFORMS Journal on Computing*, 21(2):178–192, 2009.

[10] Verena Heidrich-Meisner, Martin Lauer, Christian Igel, and Martin Riedmiller. Reinforcement learning: A friendly introduction, 2021. Accessed: 2025-05-06.

[11] Stephanie S. Harmon. Reinforcement learning: A tutorial, 2000. Accessed: 2025-05-06.

[12] Giuseppe De Pietro and Antonio Coronato. Application of reinforcement learning and deep learning in multiple-input and multiple-output (mimo) systems. *Sensors*, 22(1):309, 2022.

[13] Please consult the official documentation https://docs.python.org/3/tutorial/datastructures.html#dictionaries for details about Python dictionaries.

[14] Note there is a separate data type in Python for ordered dictionaries https://docs.python.org/3/library/collections.html#ordereddict-objects.

[15] Richard S Sutton and Andrew G Barto. *Reinforcement Learning*. Adaptive Computation and Machine Learning series. Bradford Books, Cambridge, MA, 2 edition, November 2018.

[16] https://docs.python.org/3/tutorial/datastructures.html#more-on-lists.

[17] https://docs.python.org/3/tutorial/datastructures.html#tuples-and-sequences.

[18] Majid Ghasemi, Amir Hossein Moosavi, Ibrahim Sorkhoh, Anjali Agrawal, Fadi Alzhouri, and Dariush Ebrahimi. An introduction to reinforcement learning: Fundamental concepts and practical applications, 2024.

[19] Richard S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1):9–44, August 1988.

[20] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256, 1992.

[21] Rishabh Khandelwal. Policy gradient methods in reinforcement learning, 2020. Accessed: 2025-08-12.

[22] Lilian Weng. Three fundamental flaws in common reinforcement learning algorithms (and how to fix them), 2019. Accessed: 2025-08-12.

[23] Jan Peters and Stefan Schaal. Natural actor-critic. In *Proceedings of the 25th International Conference on Machine Learning*, pages 768–775. ACM, 2008.

[24] Ruosong Wang, Simon S. Du, Lin F. Yang, and Sham M. Kakade. Is long horizon reinforcement learning more difficult than short horizon reinforcement learning? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.

[25] Junzi Zhang, Jongho Kim, Brendan O'Donoghue, and Stephen Boyd. Sample efficient reinforcement learning with REINFORCE. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 10887–10895, 2021.

[26] Andrew G. Barto, Richard S. Sutton, and Charles W. Anderson. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-13(5):834–846, 1983.

[27] Paul Brumer and Moshe Shapiro. Control of unimolecular reactions using coherent light. *Chemical Physics Letters*, 126(6):541–546, 1986.

[28] D. J. Tannor and S. A. Rice. Control of selectivity of chemical reaction via control of wave packet evolution. *The Journal of Chemical Physics*, 83(10):5013–5018, 1985.

[29] W. S. Warren, H. Rabitz, and M. Dahleh. Coherent control of quantum dynamics. *Science*, 259(5101):1581–1589, 1993.

[30] A. Assion, T. Baumert, M. Bergt, T. Brixner, B. Kiefer, V. Seyfried, M. Strehle, and G. Gerber. Control of chemical reactions by feedback-optimized phase-shaped femtosecond laser pulses. *Science*, 282(5390):919–922, 1998.

[31] J. L. Herek, W. Wohlleben, R. J. Cogdell, D. Zeidler, and M. Motzkus. Quantum control of energy flow in light-harvesting complexes. *Nature*, 417:533–535, 2002.

[32] A. P. Peirce, M. A. Dahleh, and H. Rabitz. Optimal control of quantum-mechanical systems: Existence, numerical approximation, and applications. *Physical Review A*, 37(12):4950–4964, 1988.

[33] N. Khaneja, T. Reiss, C. Kehlet, T. Schulte-Herbrüggen, and S. J. Glaser. Optimal control of coupled spin dynamics: Design of nmr pulse sequences by gradient ascent algorithms. *Journal of Magnetic Resonance*, 172(2):296–305, 2005.

[34] V. F. Krotov. *Global Methods in Optimal Control Theory.* Marcel Dekker, New York, 1996.

[35] T. Caneva, T. Calarco, and S. Montangero. Chopped random-basis quantum optimization. *Physical Review A*, 84(2):022326, 2011.

[36] D. J. Egger and F. K. Wilhelm. Adaptive hybrid optimal quantum control for imprecisely characterized systems. *Physical Review Letters*, 112(24):240503, 2014.

[37] S. Machnes, U. Sander, S. J. Glaser, P. de Fouquieres, A. Gruslys, S. Schirmer, and T. Schulte-Herbrüggen. Comparing, optimizing, and benchmarking quantum-control algorithms in a unifying programming framework. *Physical Review A*, 84(2):022305, 2011.

[38] S. J. Glaser, U. Boscain, T. Calarco, C. P. Koch, W. Köckenberger, R. Kosloff, I. Kuprov, B. Luy, S. Schirmer, T. Schulte-Herbrüggen, D. Sugny, and F. K. Wilhelm. Training schrödinger's cat: Quantum optimal control. strategic report on current status, visions and goals for research in europe. *The European Physical Journal D*, 69(12):279, 2015.

[39] A. Ansel, R. C. B. Lucherini, M. Reagor, M. J. Gullans, M. Devoret, and L. I. Glazman. Machine learning optimization of quantum control in superconducting qubits. *npj Quantum Information*, 7(94), 2021.

[40] D. Bukov, A. G. R. Day, P. Weinberg, A. Polkovnikov, and P. Mehta. Reinforcement learning in different phases of quantum control. *Physical Review X*, 8:031086, 2018.

[41] Z. Zhang, M. Sarovar, and K. B. Whaley. Deep reinforcement learning for quantum gate control. *Physical Review Letters*, 122(2):020501, 2019.

[42] A. Skolik, J. R. McClean, M. Mohseni, P. van der Smagt, and V. Dunjko. Reinforcement learning for quantum control and quantum error correction. *Quantum Science and Technology*, 7(1):015002, 2022.

[43] Y. Gao, X. Wang, N. Yu, and B. M. Wong. Harnessing deep reinforcement learning to construct time-dependent optimal fields for quantum control dynamics. *Physical Chemistry Chemical Physics*, 24(39):24208–24217, 2022.

[44] D. Koutromanos, D. Stefanatos, and E. Paspalakis. Control of qubit dynamics using reinforcement learning. *Information*, 15(5):272, 2024.

[45] J. O. Ernst, A. Chatterjee, T. Franzmeyer, and A. Kuhn. Reinforcement learning for quantum control under physical constraints. *arXiv preprint arXiv:2501.14372*, 2025.

[46] A. Author and B. Collaborators. Robust quantum control using reinforcement learning. *npj Quantum Information*, 2025. In Press.

[47] Riccardo Porotti, Antoine Essig, Benjamin Huard, and Florian Marquardt. Deep reinforcement learning for quantum state preparation with weak nonlinear measurements. *Quantum*, 6:747, 2022.

[48] Murphy Yuezhen Niu, Sergio Boixo, Vadim N. Smelyanskiy, and Hartmut Neven. Universal quantum control through deep reinforcement learning. *npj Quantum Information*, 5(1), April 2019.

[49] Thomas Fösel, Petru Tighineanu, Thomas Weiss, and Florian Marquardt. Reinforcement learning with neural networks for quantum feedback. *Quantum*, 6:747, 2022.

[50] F. Shuang and H. Rabitz. Control of quantum observables by tracking control fields. *Journal of Chemical Physics*, 121:9270–9278, 2004.

[51] D. Dong and I.R. Petersen. Quantum control theory and applications: a survey. *IET Control Theory & Applications*, 4(12):2651–2671, December 2010.

[52] Alberto Isidori. *Nonlinear Control Systems.* Springer, 1995.

[53] R. Marino and P. Tomei. Nonlinear control design: Geometric, adaptive and robust. *Prentice Hall*, 1995.

[54] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3):229–256, May 1992.

[55] John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, volume 37 of *Proceedings of Machine Learning Research*, pages 1889–1897. PMLR, 2015.

[56] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[57] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

[58] Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. *arXiv preprint arXiv:1802.09477*, 2018.

[59] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.

[60] Kenji Doya. Reinforcement learning: Computational theory and biological mechanisms. *HFSP Journal*, 1(1):30–40, May 2007.

[61] Samuel Yen-Chi Chen. An introduction to quantum reinforcement learning (qrl), 2024.

[62] Thet Htar Su, Shaswot Shresthamali, and Masaaki Kondo. Quantum framework for reinforcement learning: Integrating the markov decision process, quantum arithmetic, and trajectory search. *Physical Review A*, 111(6), June 2025.

[63] Shao hua Wang and Hai Lin. Quantum reinforcement learning, 2023.

[64] Muhammad Ismail, Mohamed Shaban, and Samuel Yen-Chi Chen. Hands-on introduction to quantum machine learning. *The International FLAIRS Conference Proceedings*, 37, May 2024.

[65] Maria Schuld and Francesco Petruccione. *Machine Learning with Quantum Computers*. Springer International Publishing, 2021.

## Appendix A: Bellman Equation Derivation

## Appendix B: Bellman Optimality Condition

In the Sec. VII, we introduced the Bellman equation. We learned that Bellman equation Eq (23) helps us in evaluating the value of a state $v_\pi(s)$ under the policy $\pi$. The Bellman optimality equation is an extension of the Bellman equation. While the standard Bellman equation expresses the value of a state $v_\pi(s)$ under a given policy $\pi$ as

$$v_\pi(s) = \sum_a \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[ r + \gamma \, v_\pi(s') \right],$$

the Bellman optimality equation describes the value of a state under the optimal policy—that is, the policy that maximizes the expected return. **Note:** We reach the optimal policy through a process known as *policy*

*improvement.* Starting with an initial policy $\pi$, we iteratively update the policy to obtain $\pi_1$, $\pi_2$, ..., until we converge to the optimal policy, which we denote by $*$ (can also denote it as $\pi_*$). As we iteratively improve the policy, the state values $v_\pi(s)$ for each state $s$ effectively increases. This process of policy improvement ensures that we converge towards the optimal policy, where the state values represent the maximum expected return.

In the optimal case, the equation is given by

$$v_*(s) = \max_a \sum_{s',r} p(s',r \mid s,a) \left[ r + \gamma \, v_*(s') \right], \quad \text{(B1)}$$

where $v_*(s)$ denotes the optimal value of state $s$.

Let us understand the above equation. For simplicity, let us consider the following situation: suppose that the agent is at $s_4$, as shown in Figure 4. Now, at this state, the action set $\mathcal{A}$ is given by

$$\mathcal{A} = \{\texttt{left}, \texttt{right}\}, \quad a \in \mathcal{A},$$

Let us now calculate the two value functions $v_\pi^{left}(s_4)$ and $v_\pi^{right}(s_4)$ corresponding to the two possible actions left and right respectively.

Using the Bellman expectation equation (note that since we define the action explicitly, the summation over the policy probabilities $\pi(a|s_4)$ disappears)):

$$v_\pi^{left}(s_4) = \sum_{s',r} p(s',r \mid s_4, left) \left[ r + \gamma v_\pi(s') \right]$$

Given the transition probabilities:

$$p(s' \neq s_3, r \mid s_4, left) = 0$$

This simplifies the value function to:

$$v_\pi^{left}(s_4) = p(s' = s_3, r = -1 \mid s_4, left) \left[ -1 + \gamma v_\pi(s_3) \right]$$
$$= p(s' = s_3, r = -1 \mid s_4, left) \times (-1 + \gamma v_\pi(s_3))$$

Similarly,
$$v_\pi^{right}(s_4) = p(s' = s_5, r = 3 \mid s_4, right) \left[ 3 + v_\pi(s_5) \right]$$

Now for the shake of simplicity let us assume the following

$$p(s' = s_5, r = 3 \mid s_4, right) = 1$$
$$p(s' = s_3, r = -1 \mid s = s_4, a = left) = 1$$

which leads to

$$
\begin{aligned}
v_\pi^{left}(s_4) &= [-1 + v_\pi(s_3)] \\
v_\pi^{right}(s_4) &= [3 + v_\pi(s_5)] \\
\text{where } v_\pi(s_3) &\leq 3 \text{ (assume } v_\pi(s_3) \leq r_{max}) \\
\text{and } v_\pi(s_5) &= 0 \text{ ( as } s_5 \text{ is optimal state)} \\
\text{Thus, } v_\pi^{left}(s_4) &\leq [-1 + 3] = 2 \\
v_\pi^{right}(s_4) &= [3 + 0] = 3
\end{aligned}
\quad \text{(B2)}
$$

According to the Bellman optimality condition, the optimal state value at $s_4$ is obtained by taking the maximum over all possible actions. That is, we have:

$$v_*(s_4) = \max\{v_\pi^{left}(s_4), v_\pi^{right}(s_4)\}.$$

From our earlier calculations, we know that:

$$v_\pi^{left}(s_4) \leq 2 \quad \text{and} \quad v_\pi^{right}(s_4) = 3.$$

Therefore, the optimal state value is:

$$v_*(s_4) = \max\{2, 3\} = 3.$$

This result indicates that the `right` action is optimal at $s_4$ (i.e., $a^* = $ `right`) because it yields the highest expected return. In summary, by applying the Bellman optimality equation, we determine that choosing the `right` action maximizes the value function at state $s_4$.

## Appendix C: Iterative Method

In this section, we will show how the Bellman equation, which was solved analytically in section VII, can be solved iteratively.

Before we show this, let us briefly recapitulate what an iterative solution method is. Consider the following equation.

$$x = f(x)$$

An iterative method aims to find the solution to this equation by starting from an initial guess and refining it step-by-step. The process is defined as:

$$x_{k+1} = f(x_k)$$

where:

- $x_k$ is the approximation of the solution at iteration $k$.

- $f(x)$ is the function that updates the current estimate.

The iterations continue until the solution converges, i.e., until the difference between consecutive estimates is below a predefined threshold (the stopping criterion):

$$|x_{k+1} - x_k| < \epsilon$$

Let $f(x) = \cos(x)$. Starting with $x_0 = 0$ and using the iteration formula $x_{\text{new}} = \cos(x_{\text{old}})$, the first three iterations are:

$$x_0 = 0,$$
$$x_1 = \cos(x_0) = \cos(0) = 1,$$
$$x_2 = \cos(x_1) = \cos(1) \approx 0.5403,$$
$$x_3 = \cos(x_2) = \cos(0.5403) \approx 0.8576.$$

Continuing the iterations until the stopping criterion is satisfied, we stop after the 10th iteration. At this point, the obtained solution is

$$x \approx 0.7314.$$

In the next section, we will present the iterative method in the case of the Bellman equation.

## Appendix D: Iterative policy evaluation

In Sec. VII, we obtained the value functions for different states by solving a system of linear equations Eq. (24) - 27. However, the same can be done using iterative method discussed in previous section. This iterative method of policy evaluation follows the equation:

$$v_\pi^{k+1}(s) = \sum_a \pi(s|a) \sum_{s',r} p(s', r \mid s, a) \left[ r + \gamma\, v_\pi^k(s') \right]$$
(D1)

Now in principle

$$\lim_{k \to \infty} v_\pi^k(s) = v_\pi(s) \quad \text{for any } s \in S \qquad \text{(D2)}$$

where $v_\pi(s)$ is the values function obtained from the solution of system of linear equation obtained under the specific policy as discussed in Sec. VII.

For instance, let's consider the task of calculating the value function $v_\pi(s_2)$ for the state $s_2$ (see Figure 4) under a given policy, denoted as $\pi$, which assigns probabilities to actions taken at state $s_2$. This policy is defined as:

$$\pi : \{\texttt{left} : p_1, \texttt{right} : p_2\}$$

where $p_1$ and $p_2$ represent the probabilities of the agent taking the `left` and `right` actions, respectively. From (D2) we have

$$\lim_{k \to \infty} v_\pi^k(s_2) = v_\pi(s_2)$$

Now the task is to calculate $v_k(s_2)$ using the iterative method. From (D1) the iterative value update equation of $v_\pi^{l+1}(s_2)$ is:

$$\forall l \in \mathbb{Z}^+ \text{ such that } l + 1 < k \in \mathbb{Z}^+$$

Therefore,

$$
\begin{aligned}
v_\pi^{l+1}(s_2) &= \sum_a \pi(s_2|a) \sum_{s',r} p(s', r \mid s_2, a) \left[ r + \gamma v_\pi^l(s') \right] \\
&= p_1 p(s_3, r \mid s_2, a) \left[ r_3 + \gamma v_\pi^l(s_3) \right] \\
&\quad + p_2 p(s_1, r \mid s_2, a) \left[ r_1 + \gamma v_\pi^l(s_1) \right] \\
&= p_1 \left[ r_3 + \gamma v_\pi^l(s_3) \right] + p_2 \left[ r_1 + \gamma v_\pi^l(s_1) \right]
\end{aligned}
$$
(D3)

where, in above derivation, we have used $p(s_3, r \mid s_2, a) = 1 = p(s_1, r \mid s_2, a)$.

In order to calculate $v_\pi^{l+1}(s_2)$, it is necessary to have the values of $v_\pi^l(s_3)$ and $v_\pi^l(s_1)$ from the previous iteration. These values can be obtained through an iterative method that updates the value functions $v_\pi^l(s)$ for all states $s \in S$ at the $l$-th iteration. The updated values are then utilized to compute the value function for the $(l+1)$-th iteration.

The initialization of the value functions can be chosen as follows:

$$v_\pi^{l=0}(s) = 0 \quad \text{for all states } s \in S.$$

.

In practice, zeros are often used because they are simple and serve as a neutral starting point, but if one has prior knowledge about the value function, it is better to choose a different initialization to potentially speed up convergence.

### Appendix E: Value iteration

Value iteration addresses the drawback of policy iteration by combining policy evaluation and improvement into a single step, performing just one update per state rather than fully evaluating the policy. This streamlined approach uses the Bellman optimality equation to ensure convergence while reducing computational cost. The equation that governs the value iteration is the following, the equation

$$v^{k+1}(s) = \max_a \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v^k(s') \right] \qquad \text{(E1)}$$

Now here we will also have the following relation

$$\lim_{k \to \infty} v^k(s) = v_*(s)$$

The value iteration equation above can also be understood in a different way: as turning the Bellman optimality equation into an update rule. Instead of performing a full policy evaluation for a fixed policy, it directly incorporates policy improvement by selecting, at each state, the action that maximizes the expected return. This streamlined approach combines evaluation and improvement into a single update step.

### Appendix F: First visit and Every visit Monte Carlo

**First visit:** When an agent encounters a state $s$ during an episode, it will record the return (sum of rewards from that point onwards) the first time it visits that state. The state's value is then updated based on that return. If the state is visited again later in the same episode, it is ignored for the purpose of updating the value estimate.

$$v(s) = \frac{1}{N(s)} \sum_{t \in \text{First Visit to } s} G_t$$

N(s) is the number of first visits to state s, $G_t$ is the return starting from time step t.

**Every visit:** Whenever the agent visits a state $s$, the value of that state is updated based on the return $G_t$ (the total reward from that point onward in the episode). The value of the state is updated multiple times within an episode if it is visited multiple times. Let the return for state $s$ at time $t$ be $G_t$, and the number of times state $s$ has been visited up to that point be $N(s)$. The value of state $s$, $V(s)$, is updated as:

$$V(s) = \frac{1}{N(s)} \sum_{t \in \text{visits to } s} G_t$$