

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/372978326>

# Reinforcement Learning: Advancements, Limitations, and Real-world Applications

Article in INTERANTIONAL JOURNAL OF SCIENTIFIC RESEARCH IN ENGINEERING AND MANAGEMENT · August 2023

DOI: 10.55041/IJSREM25118

---

CITATIONS

4

---

READS

1,758

1 author:



[Avanthikaa Srinivasan](#)

SRM Institute of Science and Technology

2 PUBLICATIONS 5 CITATIONS

SEE PROFILE

# **Reinforcement Learning: Advancements, Limitations, and Real-world Applications**

**Avanthikaa Srinivasan**

## **Abstract**

This paper aims to review the advancements, limitations, and real-world applications of RL. Additionally, it will explore the future of RL and the challenges that must be addressed to enhance its widespread applicability. By addressing these challenges, RL can be further harnessed to tackle complex real-world problems.

## **1. Introduction**

Reinforcement Learning is a subfield of machine learning that allows an agent to learn how to behave in an environment based on trial and error.

Reinforcement learning addresses the problem of how agents should learn to take actions to maximize cumulative reward through interactions with the environment. The traditional approach for reinforcement learning algorithms requires carefully chosen feature representations, which are usually hand engineered.

Reinforcement learning plays a crucial role in the field of artificial intelligence and machine learning due to its ability to handle complex decision-making tasks, adapt to changing environments and learn from its interactions. As technology advances, the importance of RL is expected to grow, paving the way for more autonomous, adaptive and intelligent systems across a wide range of applications and industries.

## **2. Background and Fundamentals of Reinforcement Learning**

### **2.1 Understanding the Principle of Reinforcement Learning**

Reinforcement Learning is a subfield of machine learning that allows an agent to learn how to interact with an environment to achieve a specific goal. The agent takes actions in the environment, the environment provides feedback which the agent uses to improve and learn its decision learning capabilities over time.

### **2.2 Key Concepts of Reinforcement Learning**

Agent: The learning entity that interacts with the environment is known as the agent. It is the learner or the decision maker in the Reinforcement Learning process.

Environment: The external context in which the agent operates and in which it interacts. It can be thought of as a dynamic system that the agent tries to understand and influence to achieve its goal.

**State:** It is a representation of the environment at a given time, containing all the relevant information that the agent needs to make decisions. It captures the current situation of the environment, including observable variables and possibly hidden or latent variables. The agent's actions are chosen based on the current state, aiming to influence future states and achieve higher rewards.

**Action:** An action represents the moves or decisions that the agent can take while interacting with the environment. Actions are based on the agent's policy which maps states to actions and guides the agent's decision-making process. The agent aims to choose actions that lead to higher rewards or desired outcomes in the environment.

**Reward:** A reward is a scalar value provided by the environment to the agent after each action. The reward serves as feedback to the agent indicating the desirability of the action taken in the given state. The agent's learning process relies on these rewards to adjust its policy and improve decision making to maximise cumulative rewards over time.

### **2.3 Markov Decision Process (MDP)**

The Markov Decision Process (MDP) is a mathematical framework used to model decision-making in situations where the outcome depends on uncertain events and the decisions made by an agent over time. MDPs are widely used in various fields, including artificial intelligence, operations research, control systems, reinforcement learning, and economics. The fundamental concepts of MDP are as follows:

**2.3.1. States (S):** MDPs involve a set of states that represent different situations or configurations in the environment. The agent operates within this environment and moves from one state to another based on its actions.

**2.3.2. Actions (A):** At each state, the agent can take a set of actions, representing the possible decisions or moves it can make.

**2.3.3. Transition Probabilities (P):** The transition probabilities define the likelihood of moving from one state to another after taking a specific action. In other words, they represent the dynamics of the environment and the uncertainty associated with state transitions.

**2.3.4. Rewards (R):** Upon taking an action in a particular state, the agent receives a numerical reward or penalty that indicates the desirability of the action in that state. The objective of the agent is to maximize the cumulative reward over time.

**2.3.5. Policy ( $\pi$ ):** A policy is a strategy that the agent follows to select actions at each state. It defines the mapping from states to actions, guiding the agent's decision-making process.

**2.3.6. Value Function (V):** The value function estimates the expected cumulative reward that the agent can achieve from a given state while following a specific policy. It is a crucial concept in MDPs as it guides the agent in making informed decisions to maximize rewards.

**2.3.7. Optimal Policy ( $\pi^*$ ):** The optimal policy is the strategy that allows the agent to obtain the maximum possible cumulative reward over time. It is the best policy among all possible policies in the MDP.

## **2.4 The Bellman Equation**

The Bellman equation essentially expresses the value of a state as the maximum expected immediate reward plus the expected discounted value of the next state, considering all possible actions that the agent can take in the current state. By iteratively applying the Bellman equation to all states in the MDP, the optimal value function and optimal policy can be determined, which helps the agent make the best decisions to maximize rewards. The Bellman equation is written as follows:

$$V(s) = \max_a (R(s,a) + \gamma \sum_{s'} P(s'|s,a) V(s'))$$

Where:

- $V(s)$  is the value of state function  $s$
- $a$  represents an action in state  $s$
- $R(s,a)$  is the immediate reward received by taking an action  $a$  in state  $s$
- $\gamma$  is the discount factor that determines the importance of future rewards compared to immediate rewards
- $P(s'|s,a)$  is the transition probability from state  $s$  to state  $s'$  after taking action  $a$

## **2.5 Popular RL Algorithms**

There are various RL algorithms that represent a subset of the vast array of methods available in the field. Each algorithm has its strengths and weaknesses, making them suitable for different types of problems and scenarios.

**2.5.1. Q-Learning:** Q-Learning is one of the most well-known and widely used algorithms in reinforcement learning. The algorithm learns an action-value function (Q-function) that estimates the expected total reward from taking a particular action in a given state. Q-Learning uses the Bellman equation to iteratively update the Q-function based on the agent's experiences in the environment. Over time, the Q-function converges to the optimal action-value function, which guides the agent to make the best decisions.

**2.5.2 Deep Q-Networks (DQNs):** Deep Q-Networks are an extension of Q-Learning that leverage deep neural networks to approximate the action-value function. DQNs use neural networks to represent the Q-function, enabling them to handle high-dimensional state spaces effectively.

**2.5.3. Policy Gradient Methods:** Policy Gradient methods are a class of model-free, policy-based algorithms that directly optimize the policy function, which maps states to actions. Unlike value-based methods, they do not rely on action-value functions. Policy Gradient methods use the gradient of an objective function to update the policy parameters, seeking to increase the expected cumulative reward. They are often more effective in dealing with continuous action spaces and have shown success in complex tasks.

**2.5.4. Proximal Policy Optimization (PPO):** Proximal Policy Optimization is a popular policy gradient method that has gained widespread attention for its stability and sample efficiency. PPO aims to update the policy parameters while ensuring that the policy does not change drastically from the previous iteration, which prevents catastrophic policy collapses. PPO has become a go-to choice for many researchers due to its strong performance and ease of implementation.

**2.5.5. Deep Deterministic Policy Gradients (DDPG):** DDPG is an actor-critic algorithm that extends the DQN architecture to continuous action spaces. It uses a deterministic policy, and the actor network learns to directly map states to continuous actions. The critic network is used to estimate the action-value function and guide the actor's updates. DDPG has been successful in various continuous control tasks.

### **3. Advancements in Reinforcement Learning**

Several recent advancements in reinforcement learning have significantly pushed the boundaries of the field. One notable trend is the development of more sample-efficient algorithms that require fewer interactions with the environment to learn effective policies. Additionally, there has been a growing interest in combining reinforcement learning with other approaches, such as unsupervised learning, imitation learning, and meta-learning, leading to promising results in learning complex tasks with limited data. Moreover, research on multi-agent reinforcement learning has advanced, enabling agents to tackle increasingly complex and interactive scenarios, including cooperative and competitive environments. While these advancements have shown impressive performance in various domains, the exploration of more scalable and interpretable RL algorithms remains an ongoing area of interest for researchers and practitioners alike.

**3.1 Deep Reinforcement Learning:** Deep Reinforcement Learning (Deep RL) is a subfield of machine learning that combines reinforcement learning (RL) with deep neural networks. In traditional RL, an agent learns to take actions in an environment to maximize a cumulative reward signal. Deep RL

enhances this process by using deep neural networks as function approximators to represent value functions or policies, enabling the agent to handle high-dimensional and complex state spaces.

**3.2 Model-Based Reinforcement Learning (Model-Based RL):** Model-Based Reinforcement Learning (Model-Based RL) involves learning an explicit model of the environment dynamics to assist in decision-making. Instead of directly interacting with the real environment, the agent simulates different scenarios using the learned model and plans its actions based on the simulations. This approach can be useful when interacting with real-world environments is costly or time-consuming.

### **3.3 Meta Reinforcement Learning (Meta RL):**

Meta Reinforcement Learning (Meta RL) deals with the problem of agents learning to learn efficiently. In other words, it focuses on developing agents that can adapt to new tasks quickly by leveraging experiences from past tasks. Meta RL algorithms aim to find representations or policies that generalize across multiple tasks, enabling more efficient learning in new, unseen tasks.

### **3.4 Multi-Agent Reinforcement Learning (Multi-Agent RL):**

Multi-Agent Reinforcement Learning (MARL) involves multiple agents interacting in a shared environment, where their actions influence each other's rewards and learning. This introduces a more complex and challenging learning scenario compared to single-agent RL, as agents must adapt to the strategies of other agents in the environment.

### **3.5 Neural networks and RL:**

Advancements in neural networks have played a crucial role in shaping the capabilities of reinforcement learning agents. Some key developments include the use of deeper architectures, attention mechanisms, and techniques for improving sample efficiency.

#### **3.5.1. Deep Architectures:**

The success of Deep Reinforcement Learning (DRL) owes much to the adoption of deep neural networks. These architectures can effectively learn complex representations from high-dimensional state spaces, enabling agents to handle real-world tasks. For example, Deep Q-Networks (DQN) utilized deep convolutional neural networks to approximate the action-value function in Atari games.

#### **3.5.2. Attention Mechanisms:**

Attention mechanisms have been instrumental in enhancing the performance of DRL agents, particularly in tasks with long sequences or complex interactions. Attention mechanisms enable agents to focus on

relevant parts of the input, which can lead to better policy decisions. Attention-based methods have shown remarkable results in tasks such as language understanding and robotic manipulation.

### 3.5.3. Sample Efficiency:

Improving sample efficiency is an essential challenge in RL, as agents often need to interact with the environment extensively to learn effective policies. Recent advancements have focused on techniques like Trust Region Policy Optimization (TRPO), Proximal Policy Optimization (PPO), and Soft Actor-Critic (SAC) to optimize policies more efficiently and with fewer samples.

### 3.5.4. Neural Network Architectures for Multi-Agent RL:

Neural networks have been extended to handle multi-agent scenarios, where agents must interact and cooperate with one another. Architectures like Multi-Agent Deep Deterministic Policy Gradients (MADDPG) and differentiable communication protocols have been proposed to tackle these challenges.

## **4. Limitations and Challenges of Reinforcement Learning**

Reinforcement Learning has seen significant advancements in neural network architectures, attention mechanisms, and sample-efficient algorithms like PPO and TRPO. However, RL still faces several limitations and challenges. Sample inefficiency remains a key issue, as agents often require a large number of interactions with the environment to learn effectively. Furthermore, RL raises ethical concerns, including bias and fairness, safety and risk management, autonomous decision-making, and resource allocation. Responsible and thoughtful deployment is crucial to ensure the safe and ethical integration of RL in society.

### **4.1. Sample Inefficiency:**

Reinforcement learning algorithms often require a large number of interactions with the environment to learn optimal policies, making them sample inefficient, especially in real-world settings where data collection can be expensive or time-consuming. Addressing this issue is crucial for the widespread adoption of RL in practical applications.

### **4.2. Exploration vs. Exploitation Trade-off:**

Balancing exploration (trying new actions to discover potentially better policies) and exploitation (leveraging already known good actions) is a fundamental challenge in RL. Ensuring that agents explore enough to discover optimal strategies while exploiting their current knowledge to maximize rewards is crucial for effective learning.

### **4.3. Reward Sparsity:**

In many real-world scenarios, the reward signals provided to the agent may be sparse, delayed, or even deceptive, making it challenging for the agent to identify the actions that lead to long-term success. Sparse rewards can hinder learning and require the development of reward shaping techniques to guide agents effectively.

### **4.4. Generalization:**

Reinforcement learning agents often struggle with generalizing their learned policies to new environments or tasks, especially when the distribution of states and rewards changes. Achieving robust and transferable policies that can adapt to different situations is an ongoing research challenge.

### **4.5 Ethical Considerations and Potential Risks of RL in the Real World:**

Reinforcement learning brings great promise, but its real-world deployment also raises ethical concerns and potential risks:

4.5.1. Bias and Fairness: RL agents learn from data, and if the data is biased, it can lead to unfair or discriminatory outcomes. Ensuring fairness and avoiding the perpetuation of existing biases is a critical concern.

4.5.2. Safety and Risk Management: In complex environments, RL agents may take actions that lead to unintended consequences or safety hazards. Ensuring the safety of RL systems and their interaction with the real world is of paramount importance.

4.5.3. Autonomous Decision Making: As RL agents become more autonomous, they may face situations that are not covered by pre-defined rules, leading to unpredictable behaviour. Ensuring accountability and responsibility for RL agents' actions is an ethical challenge.

4.5.4. Resource Allocation: In scenarios where RL is used to optimize resource allocation (e.g., in healthcare or finance), there may be ethical considerations related to the allocation of resources among different individuals or groups.

These limitations, challenges, and ethical considerations highlight the need for responsible and thoughtful deployment of reinforcement learning algorithms in real-world applications. Addressing these concerns is crucial for ensuring the safe, fair, and beneficial integration of RL in society.



## **5. Real World Applications of RL**

Despite its limitations, reinforcement learning has a wide range of real world applications due to its ability to enable agents to learn from interactions with their environment.

### **5.1. Robotics:**

Reinforcement learning has proven to be highly effective in the field of robotics, enabling autonomous systems to learn complex tasks through interactions with their environment. One prominent example is the use of reinforcement learning in robotic grasping. Instead of pre-programming specific grasping strategies, robots can learn to grasp objects of varying shapes, sizes, and materials on their own. Google's research team demonstrated this with the Dactyl robotic hand, which learned to perform a diverse range of grasping tasks through trial and error (OpenAI). These advancements have the potential to revolutionize industries like manufacturing and logistics, where robots need to adapt to ever-changing tasks and environments.

### **5.2. Finance:**

Reinforcement learning has made significant strides in the financial sector, where its ability to optimize decision-making processes and adapt to market dynamics is highly valuable. A compelling case study is trading and portfolio management. Companies like DeepMind have applied reinforcement learning to develop algorithms that autonomously learn trading strategies, optimizing investments based on market conditions (DeepMind). Additionally, reinforcement learning has been employed for personalized financial recommendations, helping individuals manage their finances better by adapting to their unique circumstances and financial goals.

### **5.3. Healthcare:**

Reinforcement learning has found compelling applications in healthcare, particularly in optimizing treatment strategies and resource allocation. For instance, in medical treatment, it can be challenging to determine the most effective dosing regimen for patients. Researchers have applied reinforcement learning to design personalized dosing policies for conditions like sepsis, aiming to improve patient outcomes while minimizing the risk of complications (Nature). Additionally, in the realm of medical imaging, reinforcement learning has been utilized to optimize image acquisition protocols, reducing radiation exposure while maintaining image quality and accuracy.

### **5.4. Gaming:**

Reinforcement learning has experienced remarkable success in gaming applications, especially in the domain of AI game agents. One standout example is AlphaGo, developed by DeepMind, which achieved unprecedented success by defeating world champion Go players. The algorithm utilized reinforcement learning to play against itself and improve its gameplay iteratively, demonstrating the potential of deep reinforcement learning in mastering complex strategic games. Additionally, reinforcement learning has

also been applied to enhance the behaviour of non-player characters (NPCs) in video games, creating more dynamic and challenging gaming experiences.

## **6. Comparison of RL with other Machine Learning Approaches:**

This section focuses on comparing and contrasting RL with other machine learning approaches, with a focus on supervised and unsupervised learning. It mainly aims to demonstrate the differences with an emphasis on objective, learning paradigm and applications. Furthermore, it discusses the advantages and disadvantages of RL.

### **6.1 Objective:**

The objective in RL is to learn a policy that maps states to actions, optimizing the agent's behaviour over time.

In supervised learning, the model learns from labelled training data, where each input is associated with a corresponding target label. The objective is to learn a mapping between inputs and outputs, enabling the model to make accurate predictions on unseen data.

Unsupervised learning involves learning patterns and structures from unlabelled data. The objective is to discover underlying relationships and representations in the data, such as clustering, dimensionality reduction, or generative modelling.

### **6.2. Learning Paradigm:**

Reinforcement learning is based on the trial-and-error learning paradigm. The agent interacts with the environment, receives feedback (rewards), and adjusts its actions to improve its performance over time. It learns from both successes and failures.

Supervised learning relies on a labelled dataset, where the model is trained on examples with known input-output pairs. The learning process involves minimizing the error between the predicted outputs and the ground truth labels.

Unsupervised learning does not use labelled data. Instead, it focuses on discovering patterns, structure, or representations in the data through techniques like clustering, dimensionality reduction, and autoencoders.

### **6.3. Applications:**

Reinforcement learning finds applications in robotics, autonomous systems, game playing, recommendation systems, finance, healthcare, and more, where agents need to learn by interacting with their environment to achieve specific goals.

**Supervised Learning:** Supervised learning is commonly used for tasks like image classification, natural language processing, sentiment analysis, and regression problems, where the model predicts target labels or values given input data.

**Unsupervised Learning:** Unsupervised learning is applied in tasks like clustering, anomaly detection, and feature learning, where the goal is to discover patterns and structures within data without labelled examples.

#### **6.4. Advantages of Reinforcement Learning (RL):**

**6.4.1. Versatility and Adaptability:** RL excels in dynamic and complex environments. Unlike supervised learning, which requires labeled data, RL agents learn directly from interactions with the environment. This adaptability allows RL to handle scenarios with changing conditions and unforeseen situations.

**6.4.2. Continuous Learning and Generalization:** RL agents can continuously learn and improve their behaviour over time. This ability is crucial for applications where the environment may evolve or when dealing with long-term tasks. RL's generalization capabilities allow it to transfer knowledge from one task to another, reducing the need for retraining from scratch.

**6.4.3. Exploration-Exploitation Tradeoff:** RL algorithms address the exploration-exploitation tradeoff, allowing agents to balance between trying new actions to discover rewards and exploiting known rewarding actions. This enables RL agents to efficiently learn optimal policies.

#### **6.5. Disadvantages of Reinforcement Learning :**

**6.5.1. Sample Inefficiency:** RL often requires a substantial number of interactions with the environment to learn effective policies, which can be computationally expensive and time-consuming. This sample inefficiency can limit RL's applicability in domains where real-world interactions are costly or dangerous.

**6.5.2. Instability and Reward Design:** Designing appropriate reward functions is a challenging aspect of RL. Incorrect or sparse reward signals can lead to instability and suboptimal policies. Tuning reward functions to guide the agent effectively is a non-trivial task and often requires domain expertise.

**6.5.3. Curse of Dimensionality:** RL's performance can degrade significantly in high-dimensional state and action spaces. The "curse of dimensionality" can make it challenging for RL agents to explore and learn efficiently, as the state space grows exponentially with the number of dimensions.

**6.5.4. Safety and Ethical Concerns:** In real-world applications, RL agents may inadvertently learn harmful or unsafe behaviours. Ensuring safety and ethical considerations in RL systems becomes crucial, especially in domains like robotics and healthcare.

## **7. Future Trends and Directions of Reinforcement Learning**

Reinforcement learning has witnessed significant advancements, but several exciting future trends are emerging, particularly in combining RL with imitation learning and transfer learning. These techniques hold promise in enhancing RL's sample efficiency, generalization, and applicability across various domains.

### **7.1. Combining RL with Imitation Learning:**

Imitation learning, also known as learning from demonstrations, involves learning a policy by observing expert behaviour. Combining RL with imitation learning can address the sample inefficiency of RL algorithms by leveraging demonstrations from experts. Researchers are exploring techniques like Behaviour Cloning, where a model is trained to imitate expert actions, and then fine-tuned using RL to improve its performance further. This approach is particularly valuable in domains where collecting RL experiences is expensive or unsafe, such as robotics and autonomous vehicles. By integrating imitation learning into RL, agents can quickly learn from expert demonstrations and fine-tune their policies through interactions with the environment.

### **7.2. Transfer Learning in Reinforcement Learning:**

Transfer learning enables agents to leverage knowledge gained from one task and apply it to related tasks. In RL, transfer learning can accelerate learning in new environments by reusing learned policies or value functions. Recent research has focused on developing methods for transferring knowledge across tasks, including meta-RL algorithms that learn how to learn efficiently. Transfer learning in RL can be especially beneficial when dealing with multi-task learning or when the agent faces a sequence of related tasks. By leveraging prior knowledge, RL agents can adapt more quickly to new environments and learn better policies.

### **7.3. Hierarchical Reinforcement Learning:**

Hierarchical RL involves learning policies at multiple levels of abstraction. Agents learn high-level policies to handle long-term objectives and low-level policies for fine-grained control. This hierarchical approach can lead to more efficient learning and better generalization across tasks. By incorporating hierarchical structures, RL agents can handle complex tasks with long horizons more effectively, making it a promising direction for real-world applications.

#### **7.4. Safe Reinforcement Learning:**

Ensuring the safety of RL agents is crucial in real-world deployments, especially in critical domains like healthcare and autonomous systems. Future trends in RL are focusing on integrating safety constraints into the learning process, leading to safe exploration and risk-sensitive policies.

Safe RL methods seek to prevent undesirable and unsafe behaviours during learning, thereby increasing the reliability and trustworthiness of RL-based systems.

### **8. Conclusion**

In conclusion, this paper explored the significant advancements, limitations, and real-world applications of reinforcement learning (RL). Over the years, RL has witnessed remarkable progress, transforming how machines learn to make decisions in dynamic environments. The integration of deep learning techniques, such as Deep Q Networks (DQNs) and policy gradients, has enabled RL agents to tackle complex tasks and achieve human-level performance in various domains.

However, despite its successes, RL still faces several challenges and limitations. Sample inefficiency remains a prominent issue, necessitating the exploration of techniques like imitation learning and transfer learning to accelerate learning and improve generalization. Moreover, the design of appropriate reward functions is often non-trivial, and RL agents may learn undesirable behaviours in safety-critical applications. Addressing these challenges is essential to unlock the full potential of RL and ensure its safe and responsible deployment in the real world.

Nevertheless, RL's real-world applications continue to grow across diverse domains. From robotics and autonomous systems to finance, healthcare, and gaming, RL has demonstrated its effectiveness in solving complex problems and optimizing decision-making processes. As research in RL advances, we can expect to witness even more innovative applications and breakthroughs, revolutionizing industries and shaping the future of AI.

### **References:**

- [1] Benchmarking Deep Reinforcement Learning for Continuous Control, Duan et al, 2016.
- [2] O'Reilly Media <https://www.oreilly.com/radar/reinforcement-learning-explained/>
- [3] <https://towardsdatascience.com/markov-decision-processes-and-bellman-equations-45234cce9d25>
- [4] Sutton, R. S., & Barto, A. G. (2018). Reinforcement Learning: An Introduction (2nd ed.). MIT Press. Chapter 3 and 4.
- [5] Puterman, M. L. (1994). Markov Decision Processes: Discrete Stochastic Dynamic Programming. John Wiley & Sons. Chapter 1.

- [6] Bellman, R. (1957). Dynamic Programming. Princeton University Press.
- [7] Watkins, C. J. C. H., & Dayan, P. (1992). Q-learning. Machine Learning, 8(3-4), 279-292.
- [8] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... & Petersen, S. (2015). Human-level control through deep reinforcement learning. Nature, 518(7540), 529-533.
- [9] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347.
- [10] Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., ... & Wierstra, D. (2016). Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971.
- [11] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. \*Advances in Neural Information Processing Systems (NeurIPS)\*, 5998-6008.
- [12] Lowe, R., Wu, Y., Tamar, A., Harb, J., Abbeel, P., & Mordatch, I. (2017). Multi-agent actor-critic for mixed cooperative-competitive environments. \*Neural Information Processing Systems (NeurIPS)\*, 6382-6393.
- [13] Kaelbling, L. P., Littman, M. L., & Cassandra, A. R. (1998). Planning and acting in partially observable stochastic domains. \*Artificial Intelligence\*, 101\*(1-2), 99-134.
- [14] Russell, S. J., & Norvig, P. (2022). Artificial intelligence: A modern approach. Pearson.
- [15] Google AI Blog: "Learning Dexterous In-Hand Manipulation" (<https://ai.googleblog.com/2018/06/scalable-deep-reinforcement-learning.html>)
- [16] OpenAI: "Dactyl - A Robotic Hand" (<https://openai.com/research/pub/dactyl>)
- [17] Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press.
- [18] Bishop, C. M. (2006). Pattern Recognition and Machine Learning. Springer.
- [19] García, J., & Fernández, F. (2015). A Comprehensive Survey on Safe Reinforcement Learning. Journal of Machine Learning Research, 16, 1437-1480.