

ARTICLE

# Mastering the Principles of Reinforcement Learning: Techniques, Applications, and Future Prospects

Firehiwot Kebede, Hailemariam Yohannes, and Getachew Desta\*

Education Strategy Center, Addis Ababa, Ethiopia.

\*Corresponding author: damiragulbada.m@mail.ru

(Received: 30 March 2023; Revised: 01 June 2023; Accepted: 05 July 2023; Published: 19 July 2023)

## Abstract

Reinforcement learning (RL) is a pivotal branch of machine learning focused on training agents to make sequences of decisions by maximizing cumulative rewards in dynamic environments. This abstract delves into the fundamental principles of RL, encompassing key techniques such as Q-learning, policy gradients, and deep reinforcement learning, which integrate neural networks to handle complex, high-dimensional tasks. RL's applications are vast and varied, extending from robotics and autonomous systems to finance, healthcare, and gaming. Notable achievements include AlphaGo's victory over human champions and the optimization of trading strategies in financial markets. The abstract also examines the challenges in RL, such as the trade-off between exploration and exploitation, scalability, and the need for substantial computational resources and data. Furthermore, the future prospects of RL are discussed, highlighting advancements in transfer learning, multi-agent systems, and the integration of RL with other machine learning paradigms to create more robust and versatile AI systems. As research progresses, mastering RL principles will be crucial for developing intelligent systems capable of adaptive, real-time decision-making, ultimately driving innovation across various sectors and transforming the landscape of artificial intelligence.

**Keywords:** Deep reinforcement learning; Exploration-exploitation; Policy gradients; Q-learning; Transfer learning; Multi-agent systems

**Abbreviations:** DQN: Deep Q Network, HER: Hindsight Experience Replay, MDP: Markov Decision Process, RL: Reinforcement learning, SAC: Soft Actor-Critic TRPO: Trust Region Policy Optimization

## 1. Introduction

Reinforcement learning is a cutting-edge field in machine learning that focuses on training intelligent agents to make optimal decisions in complex, dynamic environments. It involves an agent exploring an unknown environment through trial-and-error interactions, learning to maximize rewards by taking actions that lead to desired outcomes. The agent's goal is to discover the optimal policy, or sequence of actions, to achieve its objectives within the rules and constraints of the environment [1, 2, 3]. This comprehensive guide delves into the fundamentals of reinforcement learning, exploring key concepts such as the Markov decision process, the Bellman equation, and various reinforcement learning algorithms like Q-learning and Monte Carlo methods. It also examines the differences between reinforcement learning and supervised learning techniques, and highlights real-world applications across domains like robotics, gaming, and autonomous driving. Additionally, the guide explores advanced topics like transfer learning and its role in accelerating the training process for reinforcement learning agents [4, 5, 6, 7, 8, 9].

## 2. Fundamentals of Reinforcement Learning

Reinforcement learning (RL) is a branch of machine learning inspired by the principles of behavioral psychology. It involves an agent interacting directly with an environment, taking actions, and receiving rewards or penalties based on the outcomes of those actions. The goal of the RL agent is to learn an optimal policy, or sequence of actions, that maximizes the accumulated reward over time (Fig. 1) [10, 11, 12, 13].

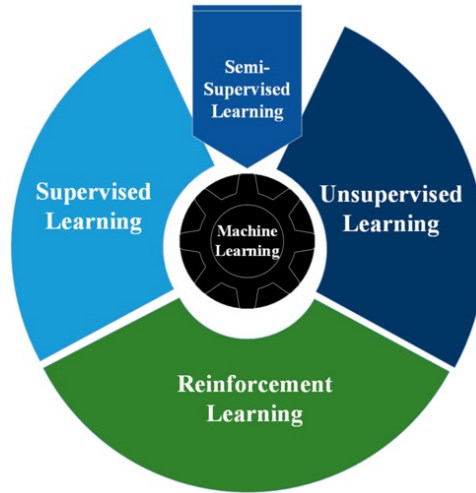


Figure 1. Machine learning branches.

The fundamental components of an RL system include:

1. **Environment:** The environment is the domain in which the agent operates, encompassing the state, actions, and rewards.
2. **Agent:** The agent is the decision-maker that interacts with the environment by taking actions and observing the resulting states and rewards.
3. **State:** The state represents the current condition or configuration of the environment.
4. **Action:** The action is the decision or behavior executed by the agent within the environment.
5. **Reward:** The reward is a numerical value that provides feedback to the agent on the desirability of the current state or action.
6. **Policy:** The policy is the strategy or function that maps states to actions, defining the agent's behavior.

The Markov Decision Process (MDP) is a mathematical framework that models the interaction between the agent and the environment over time. It consists of states, actions, rewards, and transition probabilities, which represent the likelihood of transitioning from one state to another given a specific action [14, 15, 16].

The key elements of an MDP include:

- **Value Function ( $V(s)$ ):** Represents the expected long-term reward for being in a particular state and following the optimal policy.
- **Action-Value Function ( $Q(s,a)$ ):** Represents the expected long-term reward for taking a specific action in a given state and following the optimal policy thereafter [17, 18].
- **Bellman Equation:** A fundamental equation that relates the value function to the rewards and

transition probabilities, forming the basis for dynamic programming methods to solve MDPs [19].

During the training phase, the RL agent learns to maximize the reward by repeatedly interacting with the environment and adjusting its parameters. In the inference phase, the trained RL model is deployed to perform the learned task without further parameter updates [20, 21].

### 3. Elements of Reinforcement Learning

Reinforcement learning encompasses various tasks and processes that enable an agent to learn optimal decision-making through interactions with its environment. The key elements of reinforcement learning include:

1. **Exploitation and Exploration:** The agent must strike a balance between exploiting its current knowledge to maximize rewards and exploring new actions to potentially discover better strategies.
2. **Markov Decision Processes (MDPs):** Reinforcement learning utilizes MDPs to model the decision making process, considering the current state, available actions, transition probabilities, and rewards.
3. **Sequential Decision-Making:** The agent learns through a sequential process, where each subsequent input depends on the previous decision made by the learner.
4. **Reward Maximization:** The ultimate goal of reinforcement learning is to collect as many rewards as possible by taking actions that lead to desirable outcomes.
5. **Algorithms:** Reinforcement learning employs various algorithms, such as Q-Learning, to develop solutions through step-by-step operations [22]. These algorithms often involve practical experience through commented code examples [23, 24].

**Table 1.** The core components of a reinforcement learning model

Component	Description
Policy	Determines the agent’s behavior by mapping environmental conditions to actions.
Reward	Defines the goal of the problem, providing positive or negative feedback for the agent’s actions.
Value Function	Represents the long-term attractiveness of a state based on expected future rewards.
Environment Model	Simulates the environment’s behavior, allowing the agent to predict future rewards.

The reinforcement learning process involves an agent interacting with the environment and receiving feedback in the form of rewards or punishments (Table 1). The agent is not explicitly taught what to do but must discover optimal behaviors through trial and error. Selecting the highest immediate reward may not be the best long-term strategy, as a greedy approach may not be optimal. Reinforcement learning algorithms learn from the reward/punishment feedback and adjust their behavior accordingly [25].

### 4. Reinforcement Learning Process

The reinforcement learning process involves an iterative cycle of interactions between the agent and the environment. The key steps in this process are as follows:

1. **Environment Definition:** The first step is to define the environment in which the agent will

operate. This includes specifying the state space, action space, and the rules that govern state transitions and reward calculations [26].

2. **Reward System Specification:** The reward system is a crucial component that guides the agent’s learning process. It defines the numerical rewards or penalties associated with different states and actions, enabling the agent to distinguish desirable outcomes from undesirable ones [27].
3. **Agent and Learning Algorithm Selection:** The agent is the decision-making entity that interacts with the environment. Its behavior is governed by a learning algorithm, which can be chosen from various reinforcement learning techniques such as Q-learning, SARSA, or policy gradient methods [28].
4. **Training and Validation:** During the training phase, the agent interacts with the environment, taking actions and receiving rewards or penalties based on the outcomes. The agent’s goal is to learn a policy – a mapping from states to actions – that maximizes the expected cumulative reward over time. This process involves exploration, where the agent tries out different actions to gather information, and exploitation, where the agent leverages its learned knowledge to make optimal decisions. The training process is iterative, with the agent continuously updating its policy based on the feedback received from the environment. Validation techniques, such as holdout testing or cross-validation, are employed to evaluate the agent’s performance and ensure it has learned an effective policy [29, 30].
5. **Policy Implementation:** Once the agent has learned an optimal or near-optimal policy, it can be deployed in the real-world environment or simulation to perform the desired task. The learned policy dictates the agent’s actions in response to different states encountered during the task execution [31].

Table 2. Major steps

Step	Description
Environment Definition	Specify the state space, action space, and transition rules.
Reward System Specification	Define the numerical rewards or penalties for different states and actions.
Agent and Learning Algorithm Selection	Choose the agent and the reinforcement learning algorithm it will use.
Training and Validation	Train the agent through interactions with the environment, and validate its performance.
Policy Implementation	Deploy the learned policy in the real-world or simulated environment.

From Table 2, the reinforcement learning process is iterative, with the agent continuously refining its policy through interactions with the environment until it converges to an optimal or near-optimal solution [32].

## 5. The Bellman Equation

The Bellman equation is a fundamental equation in reinforcement learning that defines the value function in terms of the rewards and transition probabilities of the Markov Decision Process (MDP). It is a recursive equation that expresses the relationship between the value of a state and the values of its successor states, along with the rewards received during the transition (Fig. 2).

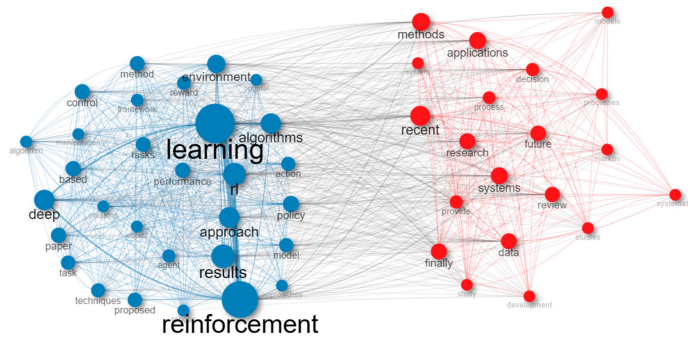


Figure 2. Co-occurrence network.

The Bellman Expectation Equation defines the value functions in terms of the immediate reward and the discounted future value [33]:

Where:

- represents the value of the current state
- is the reward received for taking an action in the current state
- is a value between 0 and 1 that determines the importance of future rewards
- is the expected value of the next state, based on the transition probabilities and the value function of the next state

This equation forms the basis for various reinforcement learning algorithms, such as Q-Learning and SARSA:

- **Q-Learning** is a model-free, off-policy algorithm based on the Bellman Equation. It uses a Q-table to store the estimated utility or quality (Q-value) of taking an action in a given state. The Q-values are iteratively updated based on the observed rewards and the maximum Q-value of the next state.
- **SARSA** (State-Action-Reward-State-Action) is a similar on-policy algorithm to Q-Learning, where the Q-values are derived from the action performed by the current policy. It updates the Q-value based on the action taken by the current policy, rather than the maximum Q-value of the next state.

Additionally, the Deep Q Network (DQN) extends Q-Learning by using neural networks to estimate the Q-value function, leveraging techniques like experience replay and target networks. This allows for more efficient learning and better generalization to complex environments [34].

## 6. Types of Reinforcement Learning

Reinforcement learning algorithms can be broadly categorized into two types: model-free RL and model-based RL.

### 6.1 Model-Free Reinforcement Learning

Model-free RL algorithms do not require a complete model of the environment. Instead, they learn directly from interactions with the environment. There are two main approaches to model-free RL [35, 36]:

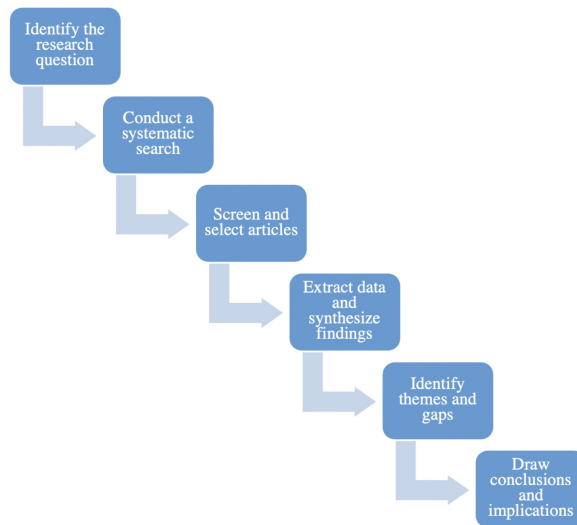
#### 1. Policy Optimization/Policy Iteration Methods:

- Policy Gradient (PG)
  - Asynchronous Advantage Actor-Critic (A3C)
  - Trust Region Policy Optimization (TRPO)
  - Proximal Policy Optimization (PPO)
2. **Q-Learning or Value-Iteration Methods:**
- Deep Q-Network (DQN)
  - C51 (Categorical DQN)
  - Quantile Regression DQN (QR-DQN)
  - Hindsight Experience Replay (HER)

Additionally, hybrid methods combine policy gradients and Q-learning, such as Deep Deterministic Policy Gradients (DDPG), Soft Actor-Critic (SAC), and Twin Delayed DDPG (TD3) [37].

## 6.2 Model-Based Reinforcement Learning

Model-based RL aims to learn or use a model of the environment to plan optimal actions. Approaches include (Fig. 1):



**Figure 3.** Research methodology.

- Learning the model
- Using the model (e.g., AlphaGo Zero)
- Hybrid methods that combine model-based and model-free techniques differs from (which requires labeled training data) and(which aims to uncover hidden structure in data). In, the agent learns from multiple trials and errors to determine the best policy or strategy to maximize rewards, making it suitable for dynamic environments where complete knowledge is not available [38].

## 7. Markov Decision Process (MDP)

A Markov Decision Process (MDP) is a mathematical framework used to model decision-making problems in dynamic, stochastic environments. It is a fundamental concept in reinforcement learning and serves as the basis for many RL algorithms. The key components of an MDP include [39, 40]:

- **States (S):** The set of possible states that the environment can be in.
- **Actions (A):** The set of actions that the agent can take in each state.
- **Transition Probabilities (P( $St+1|St, At$ )):** The probability of transitioning from one state to another, given the current state and the action taken.
- **Rewards (R( $St, At$ )):** The immediate reward received by the agent for taking an action in a specific state.

The Markov Property is a crucial assumption in MDPs, which states that the future state depends only on the current state and the action taken, and not on the previous states or actions. This property simplifies the decision-making process and enables efficient algorithms for solving MDPs (Table 3) [41, 42].

Table 3. Prime factors

Key Term	Description
Agent	The decision-maker that interacts with the environment by taking actions.
Environment	The domain in which the agent operates, encompassing states, actions, and rewards.
Policy ( $\pi$ )	The strategy or function that maps states to actions, defining the agent’s behavior.
Return	The cumulative reward received by the agent over time.
Discount Factor	A value between 0 and 1 that determines the importance of future rewards.
Value Function	Represents the expected long-term reward for being in a particular state and following the optimal policy.

The goal in an MDP is to find the optimal policy ( $\pi^*$ ) that maximizes the expected sum of discounted rewards over time. Various strategies, such as value iteration, policy iteration, SARSA, and Q-learning, can be employed to solve MDPs and find the optimal solution [43].

## 8. Reinforcement Learning Algorithms

Reinforcement learning algorithms can be broadly categorized into model-free and model-based approaches, each with its own set of techniques and algorithms [44].

### 8.1 Model-Free Reinforcement Learning Algorithms

Model-free RL algorithms do not require a complete model of the environment. They learn directly from interactions with the environment. These algorithms can be further divided into two main categories (Table 4):

#### 1. Policy Optimization/Policy Iteration Methods:

- Policy Gradient (PG)
- Asynchronous Advantage Actor-Critic (A3C)
- Trust Region Policy Optimization (TRPO)

- Proximal Policy Optimization (PPO)
2. **Q-Learning or Value-Iteration Methods:**
- Deep Q-Network (DQN)
  - C51 (Categorical DQN)
  - Quantile Regression DQN (QR-DQN)
  - Hindsight Experience Replay (HER)

Additionally, hybrid model-free methods combine policy gradients and Q-learning, such as (Table 4):

- Deep Deterministic Policy Gradients (DDPG)
- Soft Actor-Critic (SAC)
- Twin Delayed DDPG (TD3)

**Table 4.** Algorithms and role

Algorithm	Description
Deep Q-Network (DQN)	Uses neural networks to estimate Q-values, handling discrete action spaces.
Deep Deterministic Policy Gradient (DDPG)	An actor-critic algorithm that uses neural networks to approximate both the policy and value function, well-suited for continuous action spaces.
Trust Region Policy Optimization (TRPO)	An on-policy algorithm that uses neural networks to approximate the policy, ensuring conservative policy updates.
Proximal Policy Optimization (PPO)	An on-policy algorithm that uses neural networks to approximate the policy, with a clipped objective function.

Uses neural networks to estimate Q-values, handling discrete action spaces. Deep Deterministic Policy Gradient (DDPG) An actor-critic algorithm that uses neural networks to approximate both the policy and value function, well-suited for continuous action spaces. Trust Region Policy Optimization (TRPO) An on-policy algorithm that uses neural networks to approximate the policy, ensuring conservative policy updates. Proximal Policy Optimization (PPO) An on-policy algorithm that uses neural networks to approximate the policy, with a clipped objective function [45, 46, 47, 48].

## 8.2 Model-Based Reinforcement Learning Algorithms

Model-based RL algorithms aim to learn or use a model of the environment to plan optimal actions. Approaches include:

- Learning the model: Techniques like World Models, Imagination-Augmented Agents (I2A), Model-Based Priors for Model-Free RL (MBMF), and Model-Based Value Expansion (MBVE) [49].
- Using the model: Techniques employed by AlphaGo Zero, where the model is given [50].
- Hybrid methods: Combining model-based and model-free techniques. The choice of algorithm depends on factors such as the complexity of the environment, the availability of a model, and the trade-off between exploration and exploitation [51].

## 9. Reinforcement Learning vs. Supervised Learning

Reinforcement learning and supervised learning are two distinct paradigms within the field of machine learning, each with its own unique approach and applications. While supervised learning focuses on learning from labeled data, reinforcement learning emphasizes learning through interactions with an environment and receiving feedback in the form of rewards or penalties. Involves learning a generalized concept from a set of examples. It has two main tasks: regression and classification. The process involves analyzing training data, consisting of input-output pairs, to produce a generalized formula. Supervised learning algorithms, such as linear regression, logistic regression, and decision trees, aim to learn a general formula that can accurately map inputs to outputs based on the provided examples [52, 53, 54].

In contrast, does not rely on labeled data or input-output pairs. Instead, it involves an agent interacting with an environment, taking actions, and receiving rewards or penalties based on the outcomes of those actions. The goal of the agent is to learn an optimal policy, or sequence of actions, that maximizes the accumulated reward over time. This process is modeled using the Markov Decision Process (MDP), a mathematical framework that captures the dynamics of the agent-environment interaction as shown in Table 5 [55, 56].

**Table 5.** Primary differences

Aspect	Supervised Learning	Reinforcement Learning
Learning Approach	Learns from labeled data (input-output pairs)	Learns through interactions with an environment and feedback (rewards/penalties)
Goal	Learn a generalized formula to map inputs to outputs	Learn an optimal policy to maximize cumulative reward
Tasks	Regression and classification	Sequential decision-making, control mechanisms
Algorithms	Linear regression, logistic regression, decision trees	Q-learning, SARSA, policy gradients
Mathematical Framework	Analyzes training data to produce a generalized formula	MDP

While supervised learning aims to learn a general formula from the given examples, reinforcement learning focuses on controlling mechanisms and making decisions to maximize rewards in dynamic environments. Supervised learning has both input and output available during training, whereas reinforcement learning involves sequential decision-making, where the agent must learn from the consequences of its actions [57, 58, 59, 60].

## 10. Applications of Reinforcement Learning

Reinforcement learning has found widespread applications across various domains, revolutionizing fields like robotics, gaming, automation, and decision-making processes. Here are some notable applications of reinforcement learning (Fig. 4) [61, 62, 63]:

While reinforcement learning has achieved remarkable success in various domains, it also faces challenges, including the agent’s need for extensive experience, dealing with delayed rewards, and the lack of interpretability in some cases. However, ongoing advancements in deep reinforcement

learning and multi-task learning are bringing reinforcement learning closer to the realm of artificial general intelligence (AGI), further expanding its potential applications.

**Table 6.** Domains and applications

Domain	Applications
Autonomous Vehicles	<ul style="list-style-type: none"> <li>• Trajectory optimization and motion planning for self-driving cars</li> <li>• Dynamic path planning and controller optimization</li> <li>• Scenario-based learning policies for highway driving</li> <li>• Wayve.ai used deep RL to train a car to drive in a day by tackling the lane following task</li> </ul>
Robotics and Automation	<ul style="list-style-type: none"> <li>• Controlling robots to perform dangerous or repetitive tasks in industrial automation</li> <li>• Grasping and manipulating objects, including unseen ones, using techniques like QT-Opt</li> <li>• Google AI’s robots achieved a 96% success rate in grasping objects using RL</li> </ul>
Energy and Resource Management	<ul style="list-style-type: none"> <li>• DeepMind used AI agents to control Google’s data centers, leading to a 40% reduction in energy spending</li> <li>• Optimizing energy consumption and resource allocation in various industries</li> </ul>
Finance and Trading	<ul style="list-style-type: none"> <li>• Predicting stock prices and automating financial trades</li> <li>• IBM has a RL-based platform that makes financial trades and computes the reward function based on profit/loss</li> </ul>
Natural Language Processing	<ul style="list-style-type: none"> <li>• Text summarization, question answering, machine translation, and dialogue generation</li> </ul>
Healthcare	<ul style="list-style-type: none"> <li>• Determining time-dependent optimal treatment decisions for patients</li> </ul>
Engineering and Production Systems	<ul style="list-style-type: none"> <li>• Facebook’s open-source Horizon platform uses RL to optimize large-scale production systems</li> </ul>
Marketing and Advertising	<ul style="list-style-type: none"> <li>• Real-time bidding to balance competition and cooperation among advertisers</li> </ul>
Other Applications	<ul style="list-style-type: none"> <li>Game AI and game-playing (e.g., AlphaGo)</li> <li>Control theory, operations research, gaming theory, and information theory</li> <li>Synopsys uses RL in its DSO.ai solution for autonomous chip design optimization</li> <li>Traffic signal control and optimization</li> </ul>

## 11. Conclusion:

In conclusion, mastering the principles of reinforcement learning (RL) is essential for harnessing the full potential of this transformative field. RL’s core techniques, including Q-learning, policy gradients, and deep reinforcement learning, provide powerful tools for developing intelligent agents capable of complex decision-making. These methods have already demonstrated significant impact across various domains such as robotics, finance, healthcare, and gaming, showcasing the versatility and efficacy of RL in solving real-world problems. However, several challenges remain, including balancing exploration and exploitation, scalability issues, and the high demand for computational

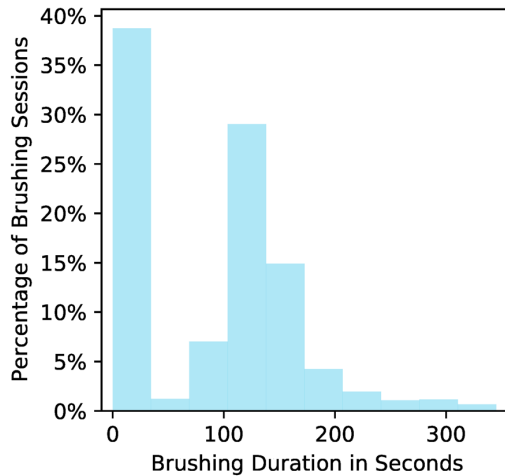


Figure 4. Histogram of brushing durations in seconds for all user.

resources and extensive data. Addressing these challenges will be crucial for the continued advancement and broader adoption of RL technologies. The future of RL looks promising, with ongoing research focusing on enhancing transfer learning, developing multi-agent systems, and integrating RL with other machine learning paradigms. These advancements are expected to lead to more robust, adaptable, and efficient AI systems capable of tackling increasingly complex tasks and environments. Ultimately, as we continue to refine and expand the principles of reinforcement learning, its applications will grow, driving innovation and shaping the future of artificial intelligence. The pursuit of mastering RL techniques will be pivotal in developing next-generation intelligent systems that can learn, adapt, and excel in a dynamic world.

## References

- [1] Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. “A brief survey of deep reinforcement learning”. In: *arXiv preprint arXiv:1708.05866* (2017).
- [2] Vaishak Belle and Ioannis Papantonis. “Principles and practice of explainable machine learning”. In: *Frontiers in big Data* 4 (2021), p. 688969.
- [3] Yang Xin, Lingshuang Kong, Zhi Liu, Yuling Chen, Yanmiao Li, Hongliang Zhu, Mingcheng Gao, Haixia Hou, and Chunhua Wang. “Machine learning and deep learning methods for cybersecurity”. In: *Ieee access* 6 (2018), pp. 35365–35381.
- [4] Yuxi Li. “Deep reinforcement learning: An overview”. In: *arXiv preprint arXiv:1701.07274* (2017).
- [5] Marco A Wiering and Martijn Van Otterlo. “Reinforcement learning”. In: *Adaptation, learning, and optimization* 12.3 (2012), p. 729.
- [6] Haoyu Yuze and He Bo. “Microbiome Engineering: Role in Treating Human Diseases”. In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 1.1 (2020), pp. 14–24.
- [7] Vincent François-Lavet, Peter Henderson, Riashat Islam, Marc G Bellemare, Joelle Pineau, et al. “An introduction to deep reinforcement learning”. In: *Foundations and Trends® in Machine Learning* 11.3-4 (2018), pp. 219–354.
- [8] Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. “Deep reinforcement learning: A brief survey”. In: *IEEE Signal Processing Magazine* 34.6 (2017), pp. 26–38.

- [9] Thanh Thi Nguyen, Ngoc Duy Nguyen, and Saeid Nahavandi. “Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications”. In: *IEEE transactions on cybernetics* 50.9 (2020), pp. 3826–3839.
- [10] Thanh Thi Nguyen, Ngoc Duy Nguyen, and Saeid Nahavandi. “Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications”. In: *IEEE transactions on cybernetics* 50.9 (2020), pp. 3826–3839.
- [11] Thanh Thi Nguyen, Ngoc Duy Nguyen, and Saeid Nahavandi. “Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications”. In: *IEEE transactions on cybernetics* 50.9 (2020), pp. 3826–3839.
- [12] Nesim Yilmaz, Tuncer Demir, Safak Kaplan, and Sevilin Demirci. “Demystifying Big Data Analytics in Cloud Computing”. In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 1.1 (2020), pp. 25–36.
- [13] Iqbal H Sarker. “Machine learning: Algorithms, real-world applications and research directions”. In: *SN computer science* 2.3 (2021), p. 160.
- [14] Yaohua Sun, Mugen Peng, Yangcheng Zhou, Yuzhe Huang, and Shiwen Mao. “Application of machine learning in wireless networks: Key techniques and open issues”. In: *IEEE Communications Surveys & Tutorials* 21.4 (2019), pp. 3072–3108.
- [15] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle-mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. “Human-level control through deep reinforcement learning”. In: *nature* 518.7540 (2015), pp. 529–533.
- [16] Peter Dayan and Nathaniel D Daw. “Decision theory, reinforcement learning, and the brain”. In: *Cognitive, Affective, & Behavioral Neuroscience* 8.4 (2008), pp. 429–453.
- [17] Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David Silver, and Koray Kavukcuoglu. “Reinforcement learning with unsupervised auxiliary tasks”. In: *arXiv preprint arXiv:1611.05397* (2016).
- [18] Jacob Oliver and William Mason. “Gene Variation: The Key to Understanding Pharmacogenomics and Drug Response Variability”. In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 1.2 (2020), pp. 97–109.
- [19] Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. “Starcraft ii: A new challenge for reinforcement learning”. In: *arXiv preprint arXiv:1708.04782* (2017).
- [20] Steven L Brunton, Bernd R Noack, and Petros Koumoutsakos. “Machine learning for fluid mechanics”. In: *Annual review of fluid mechanics* 52.1 (2020), pp. 477–508.
- [21] Fotios Zantalis, Grigorios Koulouras, Sotiris Karabetsos, and Dionisis Kandris. “A review of machine learning and IoT in smart transportation”. In: *Future Internet* 11.4 (2019), p. 94.
- [22] Chao Yu, Jiming Liu, Shamim Nemati, and Guosheng Yin. “Reinforcement learning in healthcare: A survey”. In: *ACM Computing Surveys (CSUR)* 55.1 (2021), pp. 1–36.
- [23] Jingjing Wang, Chunxiao Jiang, Haijun Zhang, Yong Ren, Kwang-Cheng Chen, and Lajos Hanzo. “Thirty years of machine learning: The road to Pareto-optimal wireless networks”. In: *IEEE Communications Surveys & Tutorials* 22.3 (2020), pp. 1472–1514.
- [24] Bauyrzhan Satipaldy, Taigan Marzhan, Ulugbek Zhenis, and Gulbadam Damira. “Geotechnology in the Age of AI: The Convergence of Geotechnical Data Analytics and Machine Learning”. In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 2.1 (2021), pp. 136–151.
- [25] Jan Peters and Stefan Schaal. “Reinforcement learning of motor skills with policy gradients”. In: *Neural networks* 21.4 (2008), pp. 682–697.
- [26] Fotios Zantalis, Grigorios Koulouras, Sotiris Karabetsos, and Dionisis Kandris. “A review of machine learning and IoT in smart transportation”. In: *Future Internet* 11.4 (2019), p. 94.

- [27] Konstantinos G Liakos, Patrizia Busato, Dimitrios Moshou, Simon Pearson, and Dionysis Bachtis. "Machine learning in agriculture: A review". In: *Sensors* 18.8 (2018), p. 2674.
- [28] Anusha Nagabandi, Ignasi Clavera, Simin Liu, Ronald S Fearing, Pieter Abbeel, Sergey Levine, and Chelsea Finn. "Learning to adapt in dynamic, real-world environments through meta-reinforcement learning". In: *arXiv preprint arXiv:1803.11347* (2018).
- [29] Mauro Birattari and Janusz Kacprzyk. *Tuning metaheuristics: a machine learning perspective*. Vol. 197. Springer, 2009.
- [30] ChoHee Kim, Donghyun Gwan, and Minh Sena Nam. "Beyond the Atmosphere: The Revolution in Hypersonic Flight". In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 2.1 (2021), pp. 152–163.
- [31] Marcus Olivecrona, Thomas Blaschke, Ola Engkvist, and Hongming Chen. "Molecular de-novo design through deep reinforcement learning". In: *Journal of cheminformatics* 9 (2017), pp. 1–14.
- [32] Jens Kober, J Andrew Bagnell, and Jan Peters. "Reinforcement learning in robotics: A survey". In: *The International Journal of Robotics Research* 32.11 (2013), pp. 1238–1274.
- [33] B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. "Deep reinforcement learning for autonomous driving: A survey". In: *IEEE Transactions on Intelligent Transportation Systems* 23.6 (2021), pp. 4909–4926.
- [34] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, et al. "A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play". In: *Science* 362.6419 (2018), pp. 1140–1144.
- [35] Chengcheng Wang, Xipeng P Tan, Shu Beng Tor, and CS Lim. "Machine learning in additive manufacturing: State-of-the-art and perspectives". In: *Additive Manufacturing* 36 (2020), p. 101538.
- [36] Ishaan Jain, Anjali Reddy, and Nila Rao. "The Widespread Environmental and Health Effects of Microplastics Pollution Worldwide". In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 2.2 (2021), pp. 224–234.
- [37] Harry Surden. "Machine learning and law: An overview". In: *Research Handbook on Big Data Law* (2021), pp. 171–184.
- [38] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, et al. "Mastering chess and shogi by self-play with a general reinforcement learning algorithm". In: *arXiv preprint arXiv:1712.01815* (2017).
- [39] Muhammad Usama, Junaid Qadir, Aunn Raza, Hunain Arif, Kok-Lim Alvin Yau, Yehia Elkhatib, Amir Hussain, and Ala Al-Fuqaha. "Unsupervised machine learning for networking: Techniques, applications and research challenges". In: *IEEE access* 7 (2019), pp. 65579–65615.
- [40] Georgios A Kaissis, Marcus R Makowski, Daniel Rückert, and Rickmer F Braren. "Secure, privacy-preserving and federated machine learning in medical imaging". In: *Nature Machine Intelligence* 2.6 (2020), pp. 305–311.
- [41] Matthew E Taylor and Peter Stone. "Transfer learning for reinforcement learning domains: A survey." In: *Journal of Machine Learning Research* 10.7 (2009).
- [42] Sarah Afiq, Maryam Fikri, Rahman Ethan, and Amsyar Isfahann. "Acknowledging the Role of Buck Converter in DC-DC Conversion". In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 3.1 (2022), pp. 287–301.
- [43] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. "Offline reinforcement learning: Tutorial, review, and perspectives on open problems". In: *arXiv preprint arXiv:2005.01643* (2020).
- [44] Shiliang Sun, Zehui Cao, Han Zhu, and Jing Zhao. "A survey of optimization methods from a machine learning perspective". In: *IEEE transactions on cybernetics* 50.8 (2019), pp. 3668–3681.

- [45] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. “Mastering the game of go without human knowledge”. In: *nature* 550.7676 (2017), pp. 354–359.
- [46] Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver, and Thore Graepel. “A unified game-theoretic approach to multiagent reinforcement learning”. In: *Advances in neural information processing systems* 30 (2017).
- [47] Seyed Sajad Mousavi, Michael Schukat, and Enda Howley. “Deep reinforcement learning: an overview”. In: *Proceedings of SAI Intelligent Systems Conference (IntelliSys) 2016: Volume 2*. Springer. 2018, pp. 426–440.
- [48] Emilia Aleksy and Veera Leevi. “Discovering the Marvels and Intricacies of Physics & Astronomy: A Journey Through Fundamental Principles and Cosmic Phenomena”. In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 3.2 (2022), pp. 342–353.
- [49] Atilim Gunes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. “Automatic differentiation in machine learning: a survey”. In: *Journal of machine learning research* 18.153 (2018), pp. 1–43.
- [50] Benjamin Sanchez-Lengeling and Alán Aspuru-Guzik. “Inverse molecular design using machine learning: Generative models for matter engineering”. In: *Science* 361.6400 (2018), pp. 360–365.
- [51] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. “Grandmaster level in StarCraft II using multi-agent reinforcement learning”. In: *nature* 575.7782 (2019), pp. 350–354.
- [52] Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. “Learning to reinforcement learn”. In: *arXiv preprint arXiv:1611.05763* (2016).
- [53] Carlo Ciliberto, Mark Herbster, Alessandro Davide Ialongo, Massimiliano Pontil, Andrea Rocchetto, Simone Severini, and Leonard Wossnig. “Quantum machine learning: a classical perspective”. In: *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 474.2209 (2018), p. 20170551.
- [54] Linnea Daniel, Sondre Robin, and Matthew Aleksander. “Future Facts: Unveiling Mental Health Issues in the Digital Age”. In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 3.2 (2022), pp. 354–365.
- [55] Jigar Patel, Sahil Shah, Priyank Thakkar, and Ketan Kotecha. “Predicting stock market index using fusion of machine learning techniques”. In: *Expert systems with applications* 42.4 (2015), pp. 2162–2172.
- [56] Keith T Butler, Daniel W Davies, Hugh Cartwright, Olexandr Isayev, and Aron Walsh. “Machine learning for molecular and materials science”. In: *Nature* 559.7715 (2018), pp. 547–555.
- [57] Benjamin Recht. “A tour of reinforcement learning: The view from continuous control”. In: *Annual Review of Control, Robotics, and Autonomous Systems* 2.1 (2019), pp. 253–279.
- [58] Jenna Burrell. “How the machine ‘thinks’: Understanding opacity in machine learning algorithms”. In: *Big data & society* 3.1 (2016), p. 2053951715622512.
- [59] Taher M Ghazal, Mohammad Kamrul Hasan, Muhammad Turki Alshurideh, Haitham M Alzoubi, Munir Ahmad, Syed Shehryar Akbar, Barween Al Kurdi, and Iman A Akour. “IoT for smart cities: Machine learning approaches in smart healthcare—A review”. In: *Future Internet* 13.8 (2021), p. 218.
- [60] Fernanda Hernández, Leonardo Sánchez, Gabriela González, and Andrés Ramírez. “Revolutionizing CMOS VLSI with Innovative Memory Design Techniques”. In: *Fusion of Multidisciplinary Research, An International Journal (FMR)* 3.2 (2022), pp. 366–379.

- [61] Sebastian Raschka, Joshua Patterson, and Corey Nolet. "Machine learning in python: Main developments and technology trends in data science, machine learning, and artificial intelligence". In: *Information* 11.4 (2020), p. 193.
- [62] Quanming Yao, Mengshuo Wang, Yuqiang Chen, Wenyuan Dai, Yu-Feng Li, Wei-Wei Tu, Qiang Yang, and Yang Yu. "Taking human out of learning applications: A survey on automated machine learning". In: *arXiv preprint arXiv:1810.13306* 31 (2018).
- [63] Emmanuel Gbenga Dada, Joseph Stephen Bassi, Haruna Chiroma, Adebayo Olusola Adetunmbi, Opeyemi Emmanuel Ajibuwa, et al. "Machine learning for email spam filtering: review, approaches and open research problems". In: *Heliyon* 5.6 (2019).