

Reinforcement Learning

Shahid Mulani¹, Rudra Bundele², Shruti Tiwari³, Neha Mahesh Shinde⁴

^{1,2,3,4} *Department of Electrical and Computer Engineering (ECE / Computer) Bharati Vidyapeeth Deemed to be University College of Engineering, Pune, Maharashtra, India*

Abstract—Reinforcement Learning (RL) is a paradigm of machine learning wherein an agent learns to make sequential decisions by interacting with an environment to maximize cumulative rewards. This paper presents a comprehensive survey of RL, covering its theoretical foundations in Markov Decision Processes, core algorithms including Q-Learning, SARSA, and Monte Carlo methods, and advanced deep RL frameworks such as Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and Actor-Critic architectures. Applications across robotics, game-playing, healthcare, and autonomous systems are reviewed, alongside open challenges including sample inefficiency, reward specification, and safe exploration. The paper concludes with directions for future research bridging RL with large foundation models.

Index Terms—Reinforcement Learning, Deep Q-Networks, Markov Decision Processes, Policy Gradient Methods, Actor-Critic, Reward Shaping, Autonomous Systems, Machine Learning.

I. INTRODUCTION

Reinforcement Learning (RL) represents a fundamental branch of machine learning that addresses the problem of decision-making under uncertainty. Unlike supervised learning, which relies on labeled datasets, or unsupervised learning, which seeks structure in unlabeled data, RL is distinguished by its interactive nature: an agent learns by receiving feedback in the form of rewards or penalties as it takes actions within an environment. This trial-and-error paradigm closely mirrors natural learning processes and has enabled breakthroughs in domains ranging from game-playing to robotics and clinical decision support.

The origins of RL trace back to early works in dynamic programming by Bellman [1] and temporal difference learning by Sutton [2]. These contributions

established the mathematical machinery that underpins modern RL. Subsequent decades witnessed the convergence of RL with deep neural networks, giving rise to deep RL methods capable of learning directly from high-dimensional sensory inputs. Landmark systems such as AlphaGo [3] and OpenAI Five demonstrated superhuman performance on tasks previously considered intractable for automated systems.

This paper is organized as follows: Section II presents the theoretical background and MDP formulation. Section III discusses classical RL algorithms. Section IV covers deep reinforcement learning methods. Section V surveys real-world applications. Section VI identifies open challenges, and Section VII concludes with future directions.

II. THEORETICAL BACKGROUND

A. Markov Decision Process

A Markov Decision Process (MDP) provides a formal framework for RL problems. Formally, an MDP is defined as a 5-tuple (S, A, P, R, γ) , where:

- S : A finite set of environment states.
- A : A finite set of possible actions.
- $P(s'|s, a)$: Transition probability function.
- $R(s, a)$: Reward function mapping state-action pairs to scalar rewards.
- $\gamma \in [0,1]$: Discount factor governing the importance of future rewards.

The agent's objective is to find an optimal policy $\pi^*(s)$ that maximizes the expected discounted return. The Markov property assumes that the future is conditionally independent of the past given the present state, enabling tractable dynamic programming solutions.

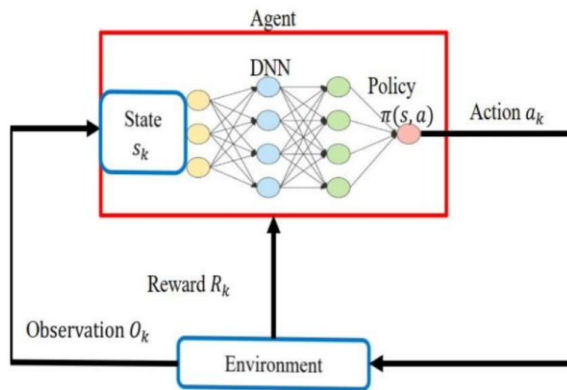
$$G_t = \sum_{k=0}^{\infty} \gamma^k R(t+k+1)$$

B. Value Functions and Bellman Equations

The state-value function $V^{\pi}(s) = E[G_t | S_t = s]$ measures the expected return from state s following policy π . The action-value function $Q^{\pi}(s, a) = E[G_t | S_t = s, A_t = a]$ measures the expected return after taking action a in state s . The Bellman optimality equation characterizes the optimal value function:

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

These equations form the basis for dynamic programming approaches and many modern RL algorithms.



III. CLASSICAL REINFORCEMENT LEARNING ALGORITHMS

A. Q-Learning

Q-Learning, introduced by Watkins and Dayan [4], is an off-policy temporal difference (TD) algorithm that directly approximates the optimal action-value function $Q^*(s, a)$.

The update rule is given by:

$$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

where α is the learning rate. Q-learning converges to the optimal policy under mild conditions, including sufficient exploration and appropriate learning rate schedules. The ϵ -greedy strategy is commonly used to balance exploration and exploitation.

B. SARSA

SARSA (State-Action-Reward-State-Action) is an on-policy TD control method. Unlike Q-Learning, SARSA updates its Q-values based on the action actually taken by the current policy, making it more conservative and better suited to environments where exploratory behavior carries cost. SARSA has demonstrated effectiveness in navigation and scheduling tasks.

C. Monte Carlo Methods

Monte Carlo (MC) methods estimate value functions by averaging complete episode returns. They do not require a model of the environment and are unbiased estimators of true value functions.

However, their reliance on complete episodes makes them unsuitable for continuing tasks and introduces high variance in estimates. MC methods are particularly useful in episodic settings such as card games and simulated environments.

IV. DEEP REINFORCEMENT LEARNING

A. Deep Q-Networks (DQN)

Mnih et al. [5] introduced Deep Q-Networks (DQN), combining Q-Learning with deep convolutional neural networks to achieve human-level performance on Atari 2600 games. Two key innovations enabled stable training: (1) experience replay, which stores transitions in a memory buffer and samples mini-batches to break temporal correlations; and (2) a target network, a periodically updated copy of the Q-network used to compute stable TD targets. Subsequent extensions including Double DQN, Dueling DQN, and Prioritized Experience Replay further improved performance.

B. Policy Gradient Methods

Policy gradient methods directly parameterize and optimize the policy $\pi_{\theta}(a|s)$. The REINFORCE algorithm [6]

estimates the policy gradient as:

$$\nabla_{\theta} J(\theta) = E_{\pi} [\nabla_{\theta} \log \pi_{\theta}(a|s) \cdot G_t]$$

While unbiased, REINFORCE suffers from high variance.

Proximal Policy Optimization (PPO) [7] and Trust Region Policy Optimization (TRPO) [8] address this

by constraining policy update magnitude, yielding stable and sample-efficient training widely adopted in applied RL systems.

C. Actor-Critic Architectures

Actor-Critic methods combine the benefits of value-based and policy-based approaches. The actor learns a policy $\pi(a|s)$ while the critic learns a value function $V_w(s)$ to reduce variance in policy gradient estimates. Asynchronous Advantage Actor-Critic (A3C) [9] enables parallel training across multiple environment instances, greatly accelerating convergence. These architectures form the backbone of state-of-the-art RL agents in continuous Model-Based Reinforcement Learning

Model-Based RL (MBRL) methods learn an explicit model of the environment dynamics $P(s'|s,a)$ and use it for planning or generating synthetic experience. Approaches such as Dyna-Q [10] and World Models [11] integrate learned models with model-free updates, substantially improving sample efficiency. MBRL is particularly valuable in data-scarce real-world domains where environment interaction is expensive.

V. APPLICATIONS OF REINFORCEMENT LEARNING

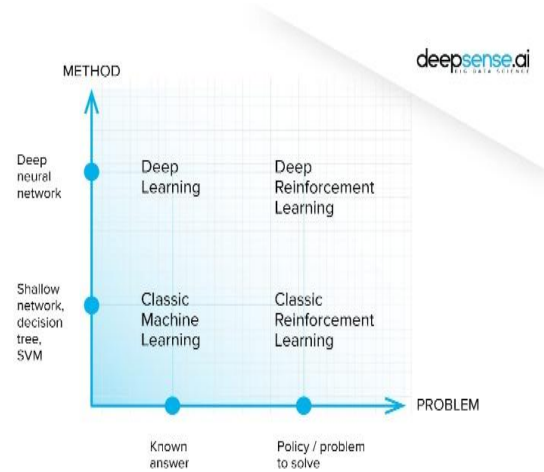
Reinforcement Learning has demonstrated transformative impact across diverse real-world domains. Table I summarizes key application areas, representative methods, and notable achievements.

Table I. Reinforcement Learning Applications by Domain

Domain	Method	Achievement
Game Playing	DQN, AlphaZero	Superhuman Atari, Chess, Go
Robotics	PPO, SAC	Dexterous manipulation
Healthcare	Actor-Critic	Treatment optimization
Autonomous Vehicles	DDPG, TD3	Path planning, control
NLP / LLMs	RLHF	Instruction following
Energy Systems	DQN variants	Grid load balancing

Reinforcement learning has particularly excelled in game environments, serving as a controlled benchmark for algorithm development. The success of AlphaZero [12] in mastering chess, shogi, and Go from self-play demonstrated that RL agents, given sufficient compute, can surpass decades of accumulated human expertise. In robotics, RL enables end-to-end policy learning for complex manipulation tasks without handcrafted controllers. In natural language processing, Reinforcement Learning from Human Feedback (RLHF) is a cornerstone technique for aligning large language models with human preferences. Model-Based Reinforcement Learning

Model-Based RL (MBRL) methods learn an explicit model of the environment dynamics $P(s'|s,a)$ and use it for planning or generating synthetic experience. Approaches such as Dyna-Q [10] and World Models [11] integrate learned models with model-free updates, substantially improving sample efficiency. MBRL is particularly valuable in data-scarce real-world domains where environment interaction is expensive.



VI. APPLICATIONS OF REINFORCEMENT LEARNING

Reinforcement Learning has demonstrated transformative impact across diverse real-world domains. Table I summarizes key application areas, representative methods, and notable achievements.

Table I. Reinforcement Learning Applications by Domain

Domain	Method	Achievement
Game Playing	DQN, AlphaZero	Superhuman Atari, Chess, Go
Robotics	PPO, SAC	Dexterous manipulation
Healthcare	Actor-Critic	Treatment optimization
Autonomous Vehicles	DDPG, TD3	Path planning, control
NLP / LLMs	RLHF	Instruction following
Energy Systems	DQN variants	Grid load balancing

Reinforcement learning has particularly excelled in game environments, serving as a controlled benchmark for algorithm development. The success of AlphaZero [12] in mastering chess, shogi, and Go from self-play demonstrated that RL agents, given sufficient compute, can surpass decades of accumulated human expertise. In robotics, RL enables end-to-end policy learning for complex manipulation tasks without handcrafted controllers. In natural language processing, Reinforcement Learning from Human Feedback (RLHF) is a cornerstone technique for aligning large language models with human preferences.

VII. CHALLENGES AND OPEN PROBLEMS

Inefficiency

One of the most critical limitations of RL is its sample inefficiency: agents often require millions of interactions to learn effective policies, a constraint prohibitive in real-world deployments. Model-based approaches, offline RL, and transfer learning are active avenues for mitigating this bottleneck.

A. Reward Function Design

Designing reward functions that accurately reflect desired behavior remains a significant engineering challenge. Poorly specified rewards can lead to reward hacking, wherein agents discover unintended shortcuts. Inverse RL and preference learning techniques aim to infer reward functions from human demonstrations or comparisons.

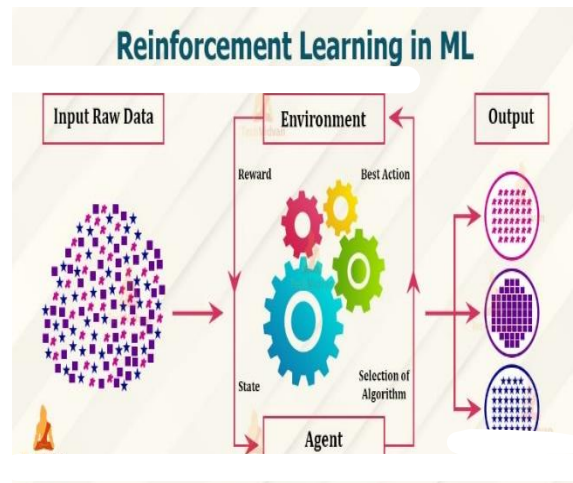
Safety and Robustness

Deploying RL agents in safety-critical systems—such as medical devices, autonomous vehicles, or industrial control—requires guarantees of safe behavior during both exploration and deployment. Constrained MDP formulations and formal verification methods are emerging approaches to address this requirement. And TransferRL agents often overfit to specific environment instances and fail to generalize to novel scenarios. Meta-reinforcement learning approaches such as MAML [13] and RL2 [14] aim to train agents that can rapidly adapt to new tasks from few interactions, a capability essential for practical deployment.

VIII. CONCLUSION AND FUTURE DIRECTIONS

This paper has presented a structured survey of reinforcement learning, covering its theoretical foundations, major algorithmic advances, and broad applicability across domains. From the Bellman equations that underpin value function estimation to deep RL architectures enabling real-world autonomous systems, RL has matured into a versatile and powerful machine learning paradigm.

Future research is expected to advance along several fronts: (1) improving sample efficiency through model-based and offline RL methods; (2) developing robust reward specification and alignment techniques; (3) establishing formal safety guarantees for deployment in high-stakes domains; and (4) enabling multi-agent and cooperative RL at scale.



ACKNOWLEDGMENT

The authors gratefully acknowledge the Department of Electrical and Computer Engineering, Bharati Vidyapeeth Deemed to be University College of Engineering, Pune, for providing the academic infrastructure and support for this research. The authors declare no conflicts of interest.

REFERENCES

- [1] R. E. Bellman, *Dynamic Programming*. Princeton, NJ: Princeton University Press, 1957.
- [2] R. S. Sutton, "Learning to predict by the methods of temporal differences," *Machine Learning*, vol. 3, no. 1, pp. 9–44, 1988.
- [3] D. Silver et al., "Mastering the game of Go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [4] C. J. C. H. Watkins and P. Dayan, "Q-Learning," *Machine Learning*, vol. 8, pp. 279–292, 1992.
- [5] V. Mnih et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, pp. 529–533, 2015.
- [6] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine Learning*, vol. 8, pp. 229–256, 1992.
- [7] J. Schulman et al., "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [8] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *Proc. ICML*, 2015, pp. 1889–1897.
- [9] V. Mnih et al., "Asynchronous methods for deep reinforcement learning," in *Proc. ICML*, 2016, pp. 1928–1937.
- [10] R. S. Sutton, "Dyna, an integrated architecture for learning, planning, and reacting," *ACM SIGART Bulletin*, vol. 2, no. 4, pp. 160–163, 1991.
- [11] D. Ha and J. Schmidhuber, "World models," *arXiv preprint arXiv:1803.10122*, 2018.
- [12] D. Silver et al., "A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play," *Science*, vol. 362, no. 6419, pp. 1140–1144, 2018.
- [13] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep

- networks," in *Proc. ICML*, 2017, pp. 1126–1135.
- [14] Y. Duan et al., "RL2: Fast reinforcement learning via slow reinforcement learning," *arXiv preprint arXiv:1611.02779*, 2016.
- [15] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. Cambridge, MA: MIT Press, 2018.