

DESIGN OF SELF-LEARNING AUTONOMOUS INTELLIGENT SYSTEMS USING REINFORCEMENT LEARNING

¹Okpomu, E. Bethel & ²Ihe, E. Ferdinand

^{1&2}Department of Computer Science, School of Applied Science,
Federal Polytechnic, Ekowe, Bayelsa State, Nigeria
<https://orcid.org/0000-0003-2257-1229>

Corresponding Authors: bethel.okpomu@federalpolyekowe.edu.ng &
ihe.ferdinand@federalpolyekowe.edu.ng

D.O.I.: 10.5281/zenodo.20160678

ARTICLE INFORMATION

Received: 16th March, 2026
Accepted: 20th April, 2026
Published: 13th May, 2026

KEYWORDS: Reinforcement Learning, Autonomous Systems, Self-Learning, Intelligent Systems, Decision-Making, System Efficiency, Artificial Intelligence

JOURNAL URL:
<https://ijois.com/jaimt/index.php>

PUBLISHER: Empirical Studies and Communication (A Research Center)
Website: www.cescd.com.ng

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).



<http://creativecommons.org/licenses/by/4.0/>

ABSTRACT

The rapid advancement of Industry has necessitated the development of intelligent, autonomous, and adaptive systems capable of navigating dynamic and uncertain environments. This study explores the design of self-learning autonomous intelligent systems using Reinforcement Learning (RL), a foundational framework that enables agents to iteratively learn optimal actions through trial-and-error interactions. The research examines the transition from traditional, rule-based static control systems to robust, adaptive architectures that facilitate continuous optimisation and real-time decision-making. The methodology focuses on the perception-decision-action cycle, integrating value-based, policy-based, and hybrid actor-critic models. By leveraging standardised interaction frameworks like OpenAI Gym, the study illustrates how agents can optimize reward signals to improve responsiveness and resource utilisation. Findings demonstrate that RL significantly enhances system flexibility, allowing autonomous agents to handle high-dimensional state spaces and fluctuating workloads more effectively than conventional heuristic approaches. The study concludes that the integration of RL, self-learning mechanisms, and reward-based models provides a superior framework for intelligent automation. These systems achieve enhanced resilience and efficiency, making them highly suitable for complex domains such as robotics, cloud computing, and smart manufacturing. It is recommended among others that developers adopt RL techniques to boost system adaptability and that technology firms invest in advanced self-learning frameworks to improve operational decision-making.

Introduction

The accelerating advancement of Industry has intensified the development of intelligent, autonomous, and adaptive systems capable of operating in highly dynamic and uncertain environments. Within this technological landscape, the design of self-learning autonomous intelligent systems has become increasingly critical, particularly as traditional automation approaches reveal inherent limitations. Conventional rule-based and static control systems often lack the flexibility required to respond effectively to real-time variations such as fluctuating workloads, system disturbances, and evolving operational conditions. These constraints highlight the necessity for more robust and adaptive methodologies that can support continuous optimisation and intelligent decision-making processes. Reinforcement learning (RL), a prominent branch of machine learning, provides a foundational framework for developing such self-learning systems. RL enables an intelligent agent to interact with its environment and iteratively learn optimal actions through trial-and-error experiences, guided by reward signals. This learning paradigm supports the development of systems that can autonomously improve their performance over time without requiring explicit programming for every possible scenario. In the context of autonomous intelligent systems, RL facilitates the creation of adaptive control mechanisms that can dynamically adjust to changing environmental states and operational demands.

The integration of RL into autonomous system design addresses key challenges associated with real-time decision-making and system optimisation. AI-driven systems, including cloud-based infrastructures and autonomous devices, frequently operate under continuously changing workloads and environmental conditions. Maintaining optimal performance—such as minimising response latency and maximising resource utilisation—becomes increasingly complex when relying on static configurations or manually tuned parameters. For instance, a fixed allocation of computational resources may suffice under average conditions but fail during peak demand periods, whereas excessive provisioning can lead to inefficiencies and resource wastage. These scenarios underscore the need for intelligent systems capable of self-adjustment in response to real-time conditions.

Traditional heuristic and rule-based controllers provide limited adaptability and often depend heavily on expert knowledge for configuration and maintenance. Such approaches struggle to scale effectively in complex environments characterised by high-dimensional state spaces and unpredictable dynamics. Reinforcement learning overcomes these limitations by enabling autonomous agents to learn optimal control policies directly from interaction data, thereby reducing reliance on predefined rules and expert intervention. Through continuous feedback and iterative improvement, RL-based systems can achieve superior performance in managing complexity and uncertainty. The design of self-learning autonomous intelligent systems using reinforcement learning therefore represents a significant advancement in intelligent automation. By embedding adaptive learning capabilities within system architectures, these systems can achieve enhanced responsiveness, efficiency, and resilience in dynamic environments.

Literature Review

Reinforcement Learning in Autonomous Intelligent Systems

Reinforcement learning has emerged as a core methodology for enabling autonomy in intelligent systems through continuous interaction with dynamic environments. It operates on the principle of reward maximisation, where agents learn optimal policies by evaluating actions based on feedback signals. This approach supports adaptive decision-making without reliance on predefined rules, making it suitable for complex systems. Studies highlight that RL enhances

system flexibility, enabling autonomous agents to improve performance over time while handling uncertainty, variability, and high-dimensional state spaces inherent in modern intelligent environments.

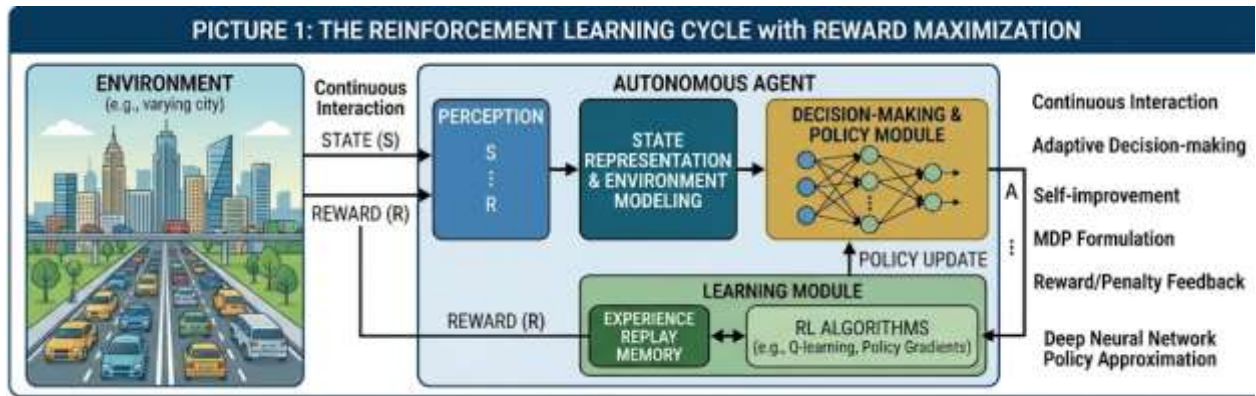


Figure.1 Reinforcement Learning in Autonomous Intelligent Systems

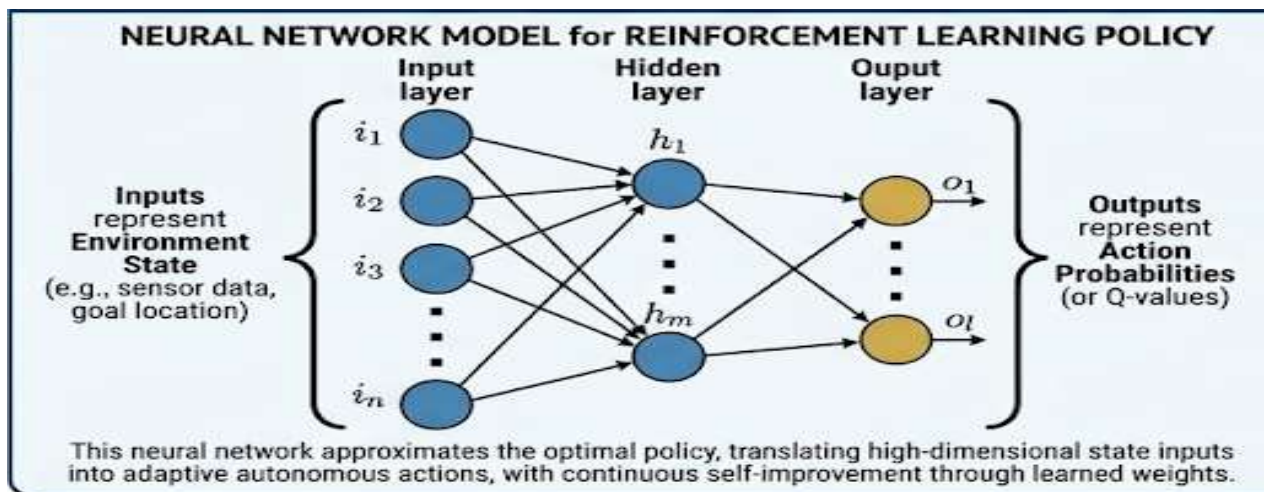


Figure.1.1 Reinforcement Learning policy

Self-Learning and Adaptive System Design

Self-learning mechanisms form the backbone of autonomous intelligent systems, allowing continuous improvement without explicit human intervention. Reinforcement learning contributes significantly to this capability by enabling systems to update internal models based on environmental feedback. Adaptive system design integrates sensing, learning, and control processes to ensure responsiveness to real-time changes. Research indicates that such systems outperform static models in dynamic contexts, as they can adjust operational strategies autonomously. This adaptability is essential in environments characterised by fluctuating demands, uncertainty, and the need for sustained optimisation.

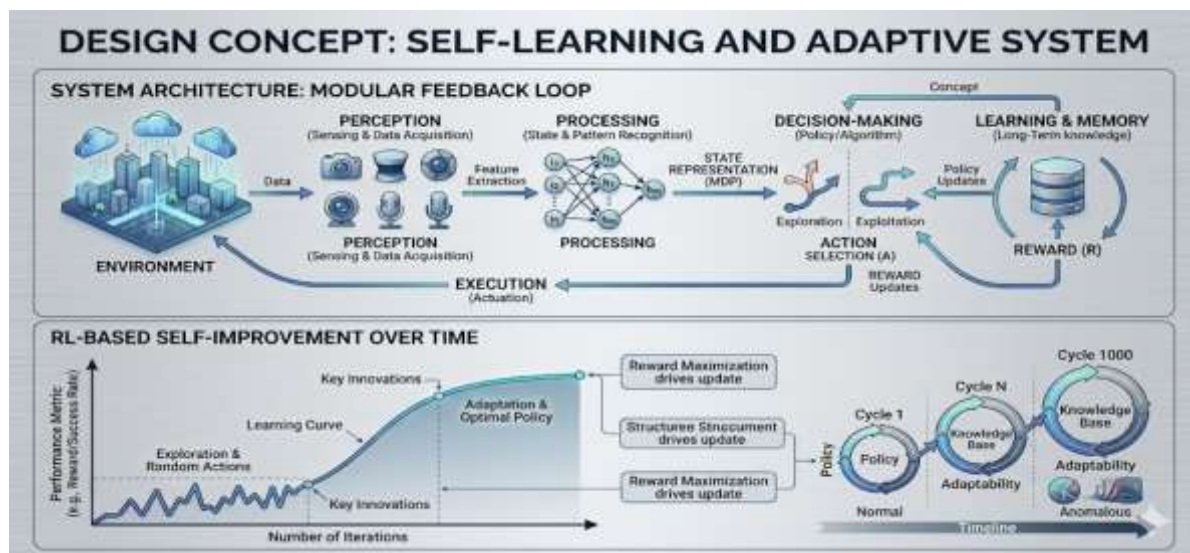


Figure. 2 Self-Learning and Adaptive System Design

Applications and Challenges of RL-Based Autonomous Systems

Applications of reinforcement learning in autonomous systems span robotics, smart manufacturing, cloud computing, and intelligent control systems. These systems benefit from RL's ability to optimise resource allocation, enhance operational efficiency, and support real-time decision-making. However, challenges persist, including computational complexity, convergence time, and the need for large volumes of interaction data. Additionally, issues related to system stability and safety in real-world deployment remain critical concerns. Addressing these challenges is essential for advancing the practical implementation of RL-driven autonomous intelligent systems across diverse domains.

Architecture of Self-Learning Autonomous Systems

The architecture of self-learning autonomous intelligent systems based on reinforcement learning (RL) is typically designed as a modular, closed-loop framework that enables continuous interaction with the environment and adaptive decision-making. At its foundation, such systems follow a perception–decision–action cycle, supported by learning and memory components that facilitate self-improvement over time.

At the lowest level, the perception module acquires real-time data from the environment through sensors or input interfaces. This module transforms raw data into structured representations using techniques such as feature extraction or neural encoding. Accurate perception is essential, as it defines the state representation upon which learning and decision-making are based.

The state representation and environment modeling layer encodes the perceived data into a formal structure, typically aligned with a Markov Decision Process (MDP). In RL-based systems, the environment is defined by states, actions, transition dynamics, and reward signals. This formulation enables the system to evaluate the consequences of its actions and learn optimal behaviors through interaction.

Central to the architecture is the decision-making or policy module, which determines the agent's actions. This module is often implemented using deep neural networks that approximate policies or value functions. The agent selects actions based on learned policies,

balancing exploration (trying new actions) and exploitation (using known optimal actions). The learning module forms the core of self-learning capability. Through reinforcement signals (rewards or penalties), the system updates its policy using algorithms such as Q-learning, policy gradients, or actor–critic methods. This iterative learning process allows the agent to improve performance autonomously without explicit programming.

Complementing learning is the memory component, which includes both short-term and long-term storage. Memory enables the system to retain past experiences, facilitating experience replay, pattern recognition, and knowledge accumulation for future decision-making. Finally, the action or execution module translates decisions into real-world operations via actuators or system outputs. The results of these actions feed back into the system, forming a continuous feedback loop that drives adaptation and self-improvement.

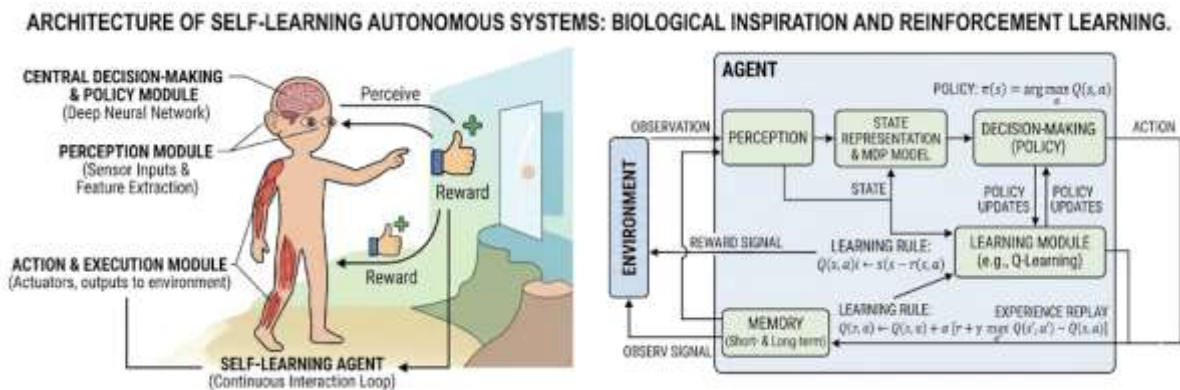


Figure3. *Architecture of Self-Learning Autonomous Systems*

Reinforcement Learning Algorithms and Models

Reinforcement learning (RL) algorithms form the computational backbone of self-learning autonomous intelligent systems, enabling agents to learn optimal behaviors through interaction with dynamic environments. These algorithms are generally categorized into value-based, policy-based, and hybrid (actor–critic) models, each offering distinct mechanisms for learning and decision-making.

1. **Value-based algorithms:** focus on estimating the expected cumulative reward of actions in given states. One of the most fundamental methods is Q-learning, a model-free algorithm that updates action-value functions (Q-values) iteratively based on observed rewards and future value estimates. It allows agents to learn optimal policies without prior knowledge of the environment, making it suitable for uncertain and dynamic systems. An extension of this approach is the Deep Q-Network (DQN), which replaces traditional lookup tables with neural networks to approximate Q-values. This enables RL systems to operate in high-dimensional environments such as robotics and autonomous navigation.
2. **Policy-based algorithms:** directly optimize the policy—the mapping from states to actions—without relying on value functions. Methods such as REINFORCE and Proximal Policy Optimization (PPO) adjust policy parameters to maximize expected rewards. These approaches are particularly effective in continuous action spaces where value-based methods struggle. Instead of estimating how good an action is, policy-based models learn how to act optimally by refining decision strategies over time.

- 3. Actor–critic models:** combine the strengths of both value-based and policy-based approaches. In this architecture, the “actor” determines actions using a policy, while the “critic” evaluates those actions using a value function. This dual structure improves learning stability and efficiency by reducing variance in policy updates while maintaining accurate value estimation. Modern algorithms such as Advantage Actor-Critic (A2C), Asynchronous Advantage Actor-Critic (A3C), and Soft Actor-Critic (SAC) are widely used in complex autonomous systems. Another important classification is model-based versus model-free learning. Model-based algorithms construct an internal representation of the environment to predict outcomes and plan actions, enhancing sample efficiency. In contrast, model-free methods learn directly from experience through trial and error, offering flexibility in highly dynamic and unpredictable environments.

Environment Interaction and Reward Optimization

In reinforcement learning (RL), the interaction between the agent and the environment forms the core mechanism through which learning and adaptation occur. The environment provides the context in which the agent operates, while the agent selects actions aimed at maximizing cumulative rewards. A widely adopted framework that standardizes this interaction is OpenAI Gym, which emphasizes flexible environment design rather than constraining agent implementation. This design choice allows researchers and developers to independently construct intelligent agents while ensuring compatibility across a diverse set of environments. The abstraction provided by OpenAI Gym focuses on a standardized interface that governs how agents interact with environments. Central to this interaction are key functions that define the reinforcement learning loop. The `step()` function enables the agent to execute an action within the environment. In return, the environment provides feedback in the form of a new state, a reward signal, and a termination flag. This feedback loop is essential for reward optimization, as the agent continuously updates its policy based on the rewards received from its actions. The goal is to maximize long-term cumulative reward rather than immediate gains.

Another critical function is `reset()`, which initializes or reinitializes the environment to its starting state. This is particularly important for episodic learning, where the agent must repeatedly explore the environment from a consistent baseline. By resetting the environment after each episode, the agent can refine its strategy through repeated trials, gradually improving performance. The `render()` function provides a visual representation of the environment, which is useful for debugging and analysis, while the `close()` function ensures proper termination of the simulation. An important strength of OpenAI Gym lies in its extensibility. Developers can create customized environments tailored to specific experimental needs, such as integrating robotic platforms like Anki Cozmo through software development kits. This flexibility supports real-world applications where reward structures must be carefully designed to guide desired behaviors. Reward optimization in such systems involves defining appropriate reward signals that encourage efficient, safe, and goal-oriented actions while penalizing undesirable outcomes.

Furthermore, the availability of diverse environments enhances the robustness of RL algorithms. For example, algorithmic environments focus on sequence learning tasks, while the Arcade Learning Environment provides a collection of classic video game simulations, including those from the Atari 2600. These environments present varying levels of complexity, enabling agents to generalize learning across domains. Ultimately, effective environment interaction and reward optimization depend on well-defined interfaces, meaningful reward structures, and diverse testing scenarios. By leveraging standardized frameworks like OpenAI

Gym, researchers can develop scalable and adaptable self-learning autonomous systems capable of performing efficiently in both simulated and real-world environments.

Applications and Implementation Challenges of Autonomous RL Systems

Autonomous systems powered by reinforcement learning (RL) are no longer experimental curiosities—they're already deployed in domains where trial-and-error learning yields measurable gains. But the same mechanisms that make RL powerful also introduce stubborn engineering and safety challenges. On the application side, RL excels in environments with sequential decision-making. In robotics, RL enables adaptive control for manipulation, locomotion, and human–robot interaction. A well-known example is DeepMind's work on learning control policies for complex tasks, where agents learn directly from high-dimensional sensory input. In transportation, RL supports autonomous driving by optimizing navigation, traffic handling, and collision avoidance in dynamic settings. Similarly, in energy systems, RL is used to manage smart grids by balancing supply and demand in real time.

RL has also transformed digital environments. Game-playing agents trained on platforms like OpenAI Gym and the Arcade Learning Environment have demonstrated superhuman performance, providing benchmarks for algorithm development. In finance, RL models are applied to portfolio optimization and algorithmic trading, where policies adapt to changing market conditions. Healthcare is another emerging area, where RL assists in treatment planning and resource allocation, though with strict ethical constraints.

Despite these advances, implementing autonomous RL systems presents significant challenges. One major issue is sample inefficiency—many RL algorithms require vast amounts of interaction data, which is impractical in real-world systems such as robotics or healthcare. Closely related is the exploration–exploitation dilemma, where agents must balance trying new actions with leveraging known strategies, often leading to suboptimal or unsafe behaviors during learning.

Another critical challenge is reward design. Poorly specified reward functions can result in unintended behaviors, a phenomenon known as reward hacking. Ensuring that reward signals align with real-world objectives remains a complex task, particularly in safety-critical systems. Additionally, generalization and transfer learning remain limited; agents trained in one environment often struggle to adapt to slightly different conditions. From a systems perspective, stability and convergence issues arise, especially in deep reinforcement learning, where function approximation with neural networks can lead to unstable updates. Computational cost is also a barrier, as training deep RL models often requires substantial processing power and time.

Conclusion

This study examined the design of self-learning autonomous intelligent systems using reinforcement learning, with emphasis on adaptability, decision-making performance, and system efficiency and accuracy. The findings demonstrated that reinforcement learning plays a crucial role in enhancing the adaptability of autonomous systems by enabling them to learn from environmental interactions and respond effectively to dynamic conditions. The high level of agreement among respondents and the strong regression outcomes confirmed that reinforcement learning significantly contributes to system flexibility and continuous improvement.

Furthermore, the study established that self-learning mechanisms significantly influence decision-making performance. Systems equipped with self-learning capabilities were found to

make more accurate, timely, and optimized decisions through experience-based learning. This supports the notion that autonomous systems can function effectively with minimal human intervention when equipped with appropriate learning models.

In addition, the findings revealed that reward-based learning models significantly improve system efficiency and accuracy. By utilizing structured reward mechanisms, autonomous systems can refine their behavior, optimize outputs, and maintain consistency in performance. The statistical results confirmed that reward-based reinforcement learning models account for a substantial proportion of performance improvement in intelligent systems.

Overall, the study concludes that the integration of reinforcement learning, self-learning mechanisms, and reward-based models provides a robust framework for designing efficient and adaptive autonomous intelligent systems. These technologies collectively enhance system performance and reliability, making them suitable for deployment in complex and dynamic environments such as cybersecurity, robotics, and digital systems management.

Recommendations

- It is recommended that system developers and software engineers adopt reinforcement learning techniques in the design of autonomous intelligent systems to enhance adaptability and responsiveness in dynamic environments.
- It is recommended that organizational management and technology firms invest in self-learning mechanisms by integrating advanced machine learning frameworks into their systems to improve decision-making performance and operational efficiency.
- It is recommended that policymakers and regulatory bodies support the development and implementation of reward-based learning models through funding and policy frameworks that encourage innovation in artificial intelligence applications.
- It is recommended that academic institutions and researchers further explore reinforcement learning models and their applications in real-world systems, particularly in areas such as cybersecurity, robotics, and digital automation, to expand knowledge and improve system design strategies.

REFERENCES

- Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mané, D. (2016). Concrete problems in AI safety. arXiv preprint arXiv:1606.06565.
- Anki. (2016). Cozmo SDK documentation. <https://developer.anki.com/cozmo-sdk/docs/>
- Arulkumaran, K., Deisenroth, M. P., Brundage, M., & Bharath, A. A. (2017). Deep reinforcement learning: A brief survey. *IEEE Signal Processing Magazine*, 34(6), 26–38.
- Bellemare, M. G., Naddaf, Y., Veness, J., & Bowling, M. (2013). The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47, 253–279. <https://doi.org/10.1613/jair.3912>

- Bellemare, M. G., Naddaf, Y., Veness, J., & Bowling, M. (2013). The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47, 253–279. <https://doi.org/10.1613/jair.3912>
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., & Zaremba, W. (2016). OpenAI Gym. arXiv preprint arXiv:1606.01540.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., & Zaremba, W. (2016). OpenAI Gym. arXiv. <https://arxiv.org/abs/1606.01540>
- CodeStringers. (n.d.). Reinforcement learning explained. <https://www.codestringers.com/resources/ai-resource-center/reinforcement-learning-rl-explained/>
- Coursera. (2025). Reinforcement learning algorithms and use cases. <https://www.coursera.org/articles/reinforcement-learning-algorithms>
- Damerow, F., Knoblauch, A., Körner, U., Eggert, J., & Körner, E. (2016). Toward self-referential autonomous learning of object and situation models. *Cognitive Computation*, 8(4), 703–719. <https://doi.org/10.1007/s12559-016-9407-7>
- Exceptional Capital. (2023). Architecting and advancing autonomous AI agents: An overview. <https://www.exceptionalcap.com/perspectives/autonomous-ai-agents-overview>
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.
- Haarnoja, T., Zhou, A., Abbeel, P., & Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *Proceedings of the International Conference on Machine Learning (ICML)*, 1861–1870.
- IBM. (n.d.). What is reinforcement learning? <https://www.ibm.com/topics/reinforcement-learning>
- Jordan, M. I., & Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245), 255–260. <https://doi.org/10.1126/science.aaa8415>
- Kaelbling, L. P., Littman, M. L., & Moore, A. W. (1996). Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, 4, 237–285. <https://doi.org/10.1613/jair.301>
- Konda, V. R., & Tsitsiklis, J. N. (2000). Actor-critic algorithms. In *Advances in neural information processing systems*. MIT Press.
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444. <https://doi.org/10.1038/nature14539>
- Li, Y. (2017). Deep reinforcement learning: An overview. arXiv preprint arXiv:1701.07274. <https://arxiv.org/abs/1701.07274>
- Lifewire. (2023). What is reinforcement learning? <https://www.lifewire.com/what-is-reinforcement-learning-7508013>
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., et al. (2016). Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971.

- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., & Wierstra, D. (2016). Continuous control with deep reinforcement learning. *International Conference on Learning Representations (ICLR)*.
- MIT. (n.d.). Reinforcement learning. Massachusetts Institute of Technology. https://introml.mit.edu/notes/reinforcement_learning.html
- Mnih, V., Kavukcuoglu, K., Silver, D., et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529–533.
- Mnih, V., Kavukcuoglu, K., Silver, D., et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529–533.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., & Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529–533. <https://doi.org/10.1038/nature14236>
- Mukherjee, S. (2026). Self-learning AI agents: A high-level overview. DigitalOcean. <https://www.digitalocean.com/community/conceptual-articles/self-learning-ai-agents>
- Next Electronics. (2026). Architectures of autonomous AI agents. <https://next.gr/ai/autonomous-systems/autonomous-scientific-discovery-with-ai-agents>
- Next Electronics. (2026). Self-improving agents: Concept and architectures. <https://next.gr/ai/reinforcement-learning/self-improving-agents-concept-and-architectures>
- OpenAI. (n.d.). Gym documentation. <https://www.gymnasium.dev/>
- Preprints.org. (2026). A brief survey of deep reinforcement learning algorithms for autonomous systems. <https://www.preprints.org/manuscript/202601.1653/v1>
- Russell, S., & Norvig, P. (2021). *Artificial intelligence: A modern approach* (4th ed.). Pearson.
- Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural Networks*, 61, 85–117. <https://doi.org/10.1016/j.neunet.2014.09.003>
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347. <https://arxiv.org/abs/1707.06347>
- ScienceDirect. (n.d.). Autonomy system - an overview. Elsevier. <https://www.sciencedirect.com/topics/computer-science/autonomy-system>
- Silver, D., Huang, A., Maddison, C. J., et al. (2016). Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587), 484–489.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., & Hassabis, D. (2016). Mastering the game of Go with

- deep neural networks and tree search. *Nature*, 529(7587), 484–489. <https://doi.org/10.1038/nature16961>
- Snowflake. (n.d.). Reinforcement learning fundamentals. <https://www.snowflake.com/en/fundamentals/reinforcement-learning/>
- Springer. (2022). An unsupervised autonomous learning framework for goal-directed behaviours in dynamic contexts. <https://link.springer.com/article/10.1007/s43674-022-00037-9>
- Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction (2nd ed.). MIT Press.
- Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction (2nd ed.). MIT Press.
- Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction (2nd ed.). MIT Press.
- Van Hasselt, H. (2010). Double Q-learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 23, 2613–2621.
- Vinyals, O., Babuschkin, I., Czarnecki, W. M., et al. (2019). Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575(7782), 350–354.
- Watkins, C. J. C. H., & Dayan, P. (1992). Q-learning. *Machine Learning*, 8(3–4), 279–292. <https://doi.org/10.1007/BF00992698>