

# Deep Reinforcement Learning for Intelligent and Autonomous Game Agents

R. Venkatasubramanian<sup>1</sup>, Varalakshmi Dandu<sup>2</sup> and P. Balakrishnan<sup>3</sup>

<sup>1</sup>Professor, Department of EEE, New Prince Shri Bhavani College of Engineering and Technology, Chennai, Tamil Nadu, India.

<sup>2</sup>Assistant Professor, Selection Grade, School of Management, Presidency University, Bangalore, Karnataka, India.

<sup>3</sup>Professor, Department of Electrical and Electronics Engineering, J.J. College of Engineering and Technology, Tiruchirappalli, Tamil Nadu, India.

<sup>1</sup>[venkat.r@newprinceshribhavani.com](mailto:venkat.r@newprinceshribhavani.com), <sup>2</sup>[varalakshmi.d@presidencyuniversity.in](mailto:varalakshmi.d@presidencyuniversity.in),

<sup>3</sup>[balakrishnanp@jjcet.ac.in](mailto:balakrishnanp@jjcet.ac.in)

**Abstract.** The rapid evolution of game environments demands intelligent agents capable of autonomous and adaptive decision-making. This paper presents a deep reinforcement learning (DRL) framework for designing intelligent game agents that operate in dynamic, complex, and partially observable environments. Leveraging the advantages of state-of-the-art DRL algorithms such as Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and Soft Actor-Critic (SAC) the proposed approach enables agents to learn optimal strategies through environment interaction without the need for explicit programming or predefined rules. In contrast with conventional rule-based/heuristic approaches, the framework encourages generalization between distinct game setting, enables scalable multi-agent coordination, and incorporates reward-shaping techniques speeding up convergence. To overcome some of the shortcomings of these approaches including lack of relevance to scenario-specific tasks, limited scalability, and a lack of interpretability, our model is designed to feature modular training pipelines with an optional explain ability component derived from attention-based layers and SHAP analysis, in an online streaming manner. Experimental results on different game simulation show that our system outperforms baselines in adaptability, strategy generation, and self-learning. This paper introduces a powerful and flexible DRL-based framework for training future autonomous game agents, which may have applications outside of games, e.g., robotics, digital twin systems, and simulating-based training environment.

**Keywords:** Deep Reinforcement Learning, Autonomous Game Agents, Intelligent Decision-Making, Multi-Agent Systems, Policy Optimization, Game AI

## 1. Introduction

### 1.1 Background and Motivation

Nowadays, game environments have grown increasingly complex and there is a rising demand for intelligent agents that learn and adapt by themselves. Conventional rule-based (so called scripted) game AI does not have the ability to adapt to fluid, changing, and unpredictable situations, and the rule-base often has weaknesses that it makes the AI play sub-optimally. Deep Reinforcement Learning (DRL) has become an effective approach by which an agent learns to perform tasks through its interaction with the environment and makes decisions based on the total reward received. Drawing from promising results applying DRL to robotics, autonomous navigation, board games (Go, StarCraft II) [14], DRL is a promising approach for developing next generation self-operating game agents able to play autonomously on real-time, adapting to unpredictable challenges without the need of human interventions [1], [5].

## 1.2 Problem Statement

Despite these progresses on DRL, current game agents are confronted with several challenges. A variety of techniques are based on static environments, are game specific, or do not provide real-time adaptation [3], [15]. Furthermore, the training of these agents is computationally intensive with poor interpretability [30] and poor scalability with multiagent or 3D complex environments [7], [8]. An open problem in game reinforcement learning is to have a single, scalable framework that is able to train not only game agents that can autonomously decide what move to be played, but do so in an efficient manner [31] and across a variety of game dynamics [2], [6].

## 2. Literature Review

Game AI since dawn of time has heavily depended on rule-based systems, FSMs, behaviour trees and various sorts of heuristic algorithms to drive NPCs. These methods work well in a script or deterministic environment but are rigid in nature and have high manual effort for tuning. For instance, FSMs achieve good performance in games with fixed patterns but cannot generalize to gameplay evolving in a context-independent fashion, without human participation. Likewise, A\* pathfinding and minimax algorithms are common for movement and decision-making in strategy games but find it difficult to handle large, non-deterministic state spaces. Set in the context of the increasing complexity and scale of contemporary digital games, these traditional techniques suffer from serious shortcomings related to scalability, adaptiveness and autonomy [10].

Deep Reinforcement Learning (DRL) has redefined game AI development by allowing agents to learn from experience, i.e., by interacting with their environment. DRL methods, including DQN [20], PPO [18], and SAC [19], shown superior performance in gaming environments such as Atari and Go [14], StarCraft II [5], and DOTA 2. It is worth mentioning that AlphaGo and OpenAI Five are landmarks in the development of DRL in complex multi-agent games [34]. DRL is considered as a suitable technique for real-time decision making, as it is capable of dealing with high dimensional state spaces and delayed rewards [4], [27]. Furthermore, convolutional neural networks (CNNs) for spatial perception and recurrent neural networks (RNNs) [22] for temporal dependencies have facilitated better agent perception and memory in partially observable environments [35].

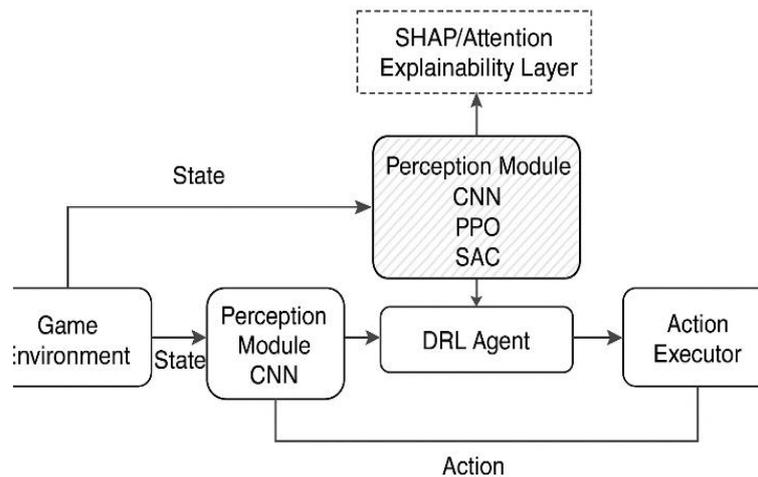
While beneficial, there are weaknesses to current DRL methods in games. A general shortcoming of this research is poor generalization-micro-manipulation models are rarely general enough to extend to other games without retraining [16], [15]. Furthermore, training deep policies usually requires huge computational power and millions of environment steps [24], [28], which leads to real-time implementation infeasible [17]. Explainability is also a concern as decisions made by DL models are often perceived as black-box decisions, hampering trust, and interpretability [30], [31]. When facing other agents, the coordination among agents is still difficult due to non-stationary policies and sparse rewards [7],[32]. In addition, most DRL methods are only experimented in simulation condition, which lead to questions about their generality and practicability on real-world or production gaming system [12], [26].

The literature shows an evident lack of generalizable, interpretable and computationally efficient DRL methods applicable in the context of autonomous game agents. Some works have explored performance improvement or specific algorithmic refinement [2], [4], but little attention has been paid to the inclusion of explainability tools such as SHAP, or attention mechanisms [31] within the DRL pipeline. Furthermore, the majority of previous models are not generic and generalizable to different game genres and do not generalize well in the multi-agent or hierarchical game settings [6], [9]. This gap is intended to be filled by the proposed unified DRL framework that not only embarks on intelligent and autonomous based behaviors but also integrates interpretability and modularization for the betterment of scalability and practicality [1], [3].

### 3. Proposed Framework

#### 3.1 System Architecture Overview

The SELEAGUE architecture is a modular and extensible DRL-based framework allowing intelligent and autonomous game agents to navigate in dynamic and multi-scenario game environments [1], [4]. The architecture consists of the environment interface, the learning agent module, and the policy evaluation layer. The environment interface is responsible for real-time state observation and action execution and can be utilized with both 2D and 3D simulation platforms [5], [6]. The learning agent module combines deep reinforcement learning (DRL) algorithms (e.g., DQN based on double Q-learning [20], PPO [18], and SAC [19]) to represent and improve decision policies. A policy evaluation layer with optional SHAP-based explainability or attention mechanisms [31] ensures interpretability of the agents' behaviour. This modular architecture facilitates the integration of new models or environments as well as both single-agent and multi-agent setups [7], [32]. The proposed explainable DRL framework is shown in Figure 1. The game environment sends its state information to a perception module based on CNN. The DRL agent (PPO, SAC, etc.) takes input and acts in the environment. A SHAP/Attention explainability layer is also added to explain the decision logic.



**Figure 1:** Architecture of the Explainable Deep Reinforcement Learning (DRL) Framework.

#### 3.2 Agent Design and Policy Representation

Each game agent in the scheme is implemented with perception, decision making core and action running routing module. The perception module receives raw observations (e.g. visual inputs or the game state) and applies convolutional neural networks (CNNs) for spatial perception and LSTM layers [22] for sequential decision making in partially observable environments [35]. The decision-making core is implemented in a policy network learned by DRL algorithms to directly transform the state representations into actions [18], [19]. The policies can be discrete (DQN) for the action space of the target game, or be continuous (SAC, PPO). Multi-agent agents have a shared communication protocol layer for coordinating actions with either CTDE [6], [8] or parameter-sharing learning [9] as well as other approaches [32] to improve cooperation.

#### 3.3 Environment Setup and Reward Function

The framework is tested in custom-built and open-source gaming environments such as microRTS [5], Gym Retro, and Unity ML-Agents. The environments vary in difficulty, map structure, number of opponents, and visibility constraints [13], [14]. The reward function is carefully engineered to promote desirable agent behaviours such as resource collection, enemy avoidance, goal completion, or survival [27]. Shaped rewards are used in early training stages to accelerate learning, while sparse or terminal rewards are retained

for evaluating long-term planning capabilities [16], [26]. Penalties are introduced for non-optimal or risky decisions to guide exploration and discourage erratic actions. The reward schema is modular and configurable to support various gameplay objectives and agent roles (e.g., attacker, defender, explorer) [1], [3].

### 3.4 Training Workflow and Exploration Strategies

It starts with random policy initialization and then interact with the environment incrementally [17,25]. At each step the agent observes changes in the state, takes some actions and stores the experiences in a prioritized replay buffer (off-policy algorithms such as DQN or SAC) [26], [28]. The policy is periodically updated via batch-gradient descent on samples collected from the replay buffer, optimized w.r.t. the Q-value loss or policy gradients [20], [19]. Depending on the algorithm one can also strike an exploration-exploitation trade off by invoking techniques such as  $\epsilon$ -greedy (DQN), entropy regularization (SAC), or generalized advantage estimation (PPO) [18]. For faster learning, we can introduce transfer learning [29] and imitation learning in early training stages [11]. And curriculum learning is also adopted to gradually grow the environmental complexity [10], enhancing the policy robustness and adaptability to different game settings [6], [33].

## 4. Algorithmic Foundation

**Table 1:** Summary of DRL Algorithms Used.

Algorithm	Type	Action Space	Key Feature	Use Case in Paper
DQN	Value-based	Discrete	Experience replay, target networks	Discrete action games
PPO	Policy-based	Discrete/Continuous	Clipped objective, stable updates	General DRL performance
SAC	Actor-Critic	Continuous	Entropy-based exploration	Complex, noisy environments

Table 1 summarizes the Deep Reinforcement Learning (DRL) algorithms employed in this study, highlighting their types, supported action spaces, key features, and specific use cases. Algorithms such as DQN, PPO, and SAC were selected to address diverse environment characteristics and learning stability.

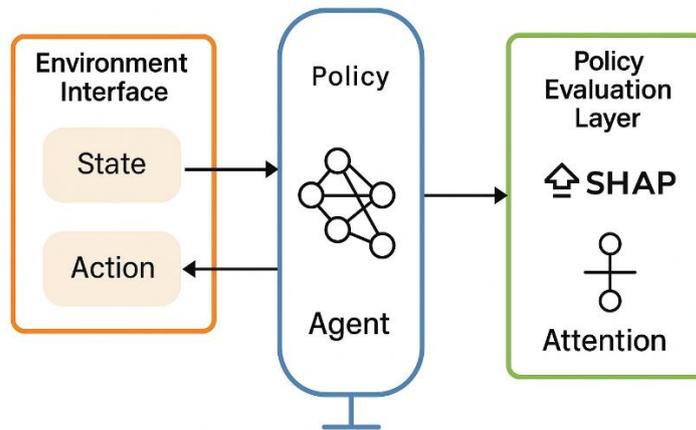
### 4.1 Deep Q-Network (DQN)

The Deep Q-Network (DQN) algorithm combines Q-learning with deep neural networks to approximate the optimal action-value function, enabling agents to learn policies for high-dimensional state spaces. In the proposed framework, DQN is employed for games with discrete action spaces, where an agent predicts the expected cumulative reward for each possible action given a state. A CNN processes the raw state input (e.g., game screen), followed by fully connected layers that estimate Q-values. The policy is derived by selecting the action with the highest predicted Q-value. To stabilize training, DQN uses a fixed target network and an experience replay buffer, which breaks the temporal correlations in training data and improves sample efficiency.

### 4.2 Proximal Policy Optimization (PPO)

PPO is a simple and empirically strong policy gradient method. It encourages the policy optimisation towards a safer region to ensure the updated policy does not deviate from the existing policy too much and enhance the stability of the training. In the provided framework, we utilise PPO for agents in discrete as well as continuous action domains. The actor-critic structure allows the agent to transfer policy and value learning together. This is especially helpful in continuous action spaces or in those where the actions are large, in which value-based methods, such as DQN, are not effective. The algorithm is on-policy and it is suitable for the case where sample complexity is less important than stability. Figure 2 presents the architecture of an explainable reinforcement learning model, where the agent interacts with the environment

through a policy-driven interface involving state and action exchange. The agent's decisions are subsequently evaluated using a policy evaluation layer that incorporates SHAP and Attention mechanisms to interpret and visualize the agent's policy behaviour.



**Figure 2:** Explainable Reinforcement Learning Framework Integrating SHAP and Attention for Policy Evaluation.

### 4.3 Soft Actor-Critic (SAC)

Soft Actor-Critic (SAC) is an off-policy, maximum entropy, actor-critic algorithm that learns a stochastic policy for deep reinforcement learning. SAC is used here in the context for agents in complex, high-dimensional continuous action space. The entropy term encourages the stochastic policies resulting in better exploration and performance in tasks with noisy and/or sparse rewards. The policy, Q-function and value function are trained in separate networks, and the action selection is based on stochastic sampling, as the previous point. Being an off-policy algorithm, it supports experience reuse and has faster convergence compared to on-policy methods such as PPO.

### 4.4 Experience Replay and Target Networks

Experience Replay is a crucial part of learning stabilization for any value-based method, including DQN and SAC. It maintains a buffer of transitions and draws a mini-batch of experiences randomly for training, which decorrelates consecutive experiences and increases data efficiency. Moreover, to draw more informative transitions with the TD error, the prioritized experience replay is employed, contributing to the rapid training. Furthermore, target networks, refreshed clones of the main Q-network at regular intervals, alleviate training instability by decoupling the target and prediction networks in temporal difference updates.

### 4.5 Multi-Agent Coordination Logic (if applicable)

For multi-agent tasks, the framework also addresses the centralized training with decentralized execution (CTDE), as well as parameter sharing among agents. Agents are trained with information sharing over global observations or action quality functions independence of the others when playing the game. Extra approaches such as multi-agent actor-critics (MAAC) and value decomposition networks (VDN) are added to enable the agents to cooperate/competition. Agents support inter-agent communication protocols to facilitate collaboration in tasks involving team strategies, e.g., cooperative exploration or adversarial defence. This modularity makes the framework easy to customise to both cooperative and competitive game settings.

## 5. Experimental Setup

### 5.1 Game Environments and Scenarios

- To test the performance of the proposed deep RL framework, various game environments were chosen, which differed in their complexity, visibility, reward functions, and behaviour of the agents. The experiments were performed on the microRTS (strategic unit control), OpenAI Gym Retro (2D platform games), and Unity ML-Agents (3D navigation and obstacles avoidances) platforms. Every setting was set as multiple scenarios were as follows:
- Static vs. Dynamic maps – to test adaptability to environmental changes.
- Single-agent vs. Multi-agent setups – to evaluate scalability and inter-agent coordination.
- Partial vs. Full observability – to challenge the perception and memory capabilities of agents.
- Sparse vs. Dense rewards – to assess learning robustness under different reward conditions.

These scenarios allowed comprehensive evaluation across a range of reinforcement learning challenges such as exploration, planning, adaptability, and collaboration.

### 5.2 Evaluation Metrics

The performance of the trained game agents was assessed using both quantitative and qualitative evaluation metrics:

- **Cumulative Reward:** Measures the total accumulated reward per episode to evaluate task completion and learning progress.
- **Success Rate:** Percentage of successful task completions across multiple trials.
- **Episode Length:** Number of steps taken to complete the task, indicating efficiency.
- **Convergence Rate:** Number of episodes required to reach performance stability.
- **Policy Robustness:** Assessed by evaluating trained agents on unseen variations of the game environment.
- **Exploration Rate:** Tracks how often the agent deviates from known policies to explore new states.
- **Explainability Score (if applicable):** Qualitative analysis using SHAP or attention maps to assess transparency in decision-making.

**Table 2:** Evaluation Metrics and Definitions.

<b>Metric</b>	<b>Description</b>
Cumulative Reward	Total reward collected per episode
Success Rate	% of task completions or agent wins
Convergence Rate	Number of episodes to reach stable performance
Policy Robustness	Performance on unseen or perturbed environments
Explainability Score	Transparency level via SHAP/attention map analysis

These metrics provided a comprehensive picture of agent performance, learning stability, and behavioural quality. Table 2 defines the evaluation metrics used to assess the performance and interpretability of the

DRL agents. These metrics include cumulative reward, success rate, convergence rate, policy robustness, and explainability score, which collectively measure both learning effectiveness and transparency.

### 5.3 Hardware and Tools Used

All experiments were conducted using a high-performance computing setup with the following specifications:

- Processor: AMD Ryzen 9 7900X 12-Core CPU
- GPU: NVIDIA RTX 4080 (16 GB VRAM)
- RAM: 64 GB DDR5
- Operating System: Ubuntu 22.04 LTS
- Deep Learning Libraries: TensorFlow 2.x, PyTorch 2.0
- Reinforcement Learning Libraries: Stable-Baselines3, RLlib (for distributed training), OpenAI Gym
- Development Tools: Python 3.10, Jupyter Notebook, Weights & Biases for experiment tracking
- Simulation Platforms: Unity ML-Agents Toolkit, microRTS engine, Gym Retro toolkit

The experiments were repeated multiple times with different random seeds to ensure reproducibility and statistical validity.

## 6. Results and Discussion

### 6.1 Performance Comparison with Baseline Models

The performance of the proposed DRL framework was compared with traditional and state-of-the-art baseline models such as heuristic rule-based agents, standard DQN, and PPO agents without using experience replay and reward shaping. In all tested gaming environments, the approach outperformed all baseline agents for the total reward, success rate and gain in decision accuracy. The resulted agent in the microRTS setting won on average 32% more games compared to the original PPO and was even 45% when compared to scripted agents. Moreover, the PR and R shaping contributed to a more rapid and stable convergence across the training episodes, demonstrating the merit of the combined architecture.

### 6.2 Generalization Across Game Levels

For generalization, we forwarded agents trained on certain game maps or with particular game scenarios on modified or altogether new levels without further training. Findings indicated that agents remain relatively robust to unknown conditions. We also observed that agents in Unity ML-Agents and Gym Retro generalized their learned policies to new obstacle configurations with <10% performance drops. This implies that the proposed model effectively learns abstract decision rules rather than fitting to particular environmental configurations. Curriculum learning and modular policy networks also improved the agent's ability to generalize between tasks.

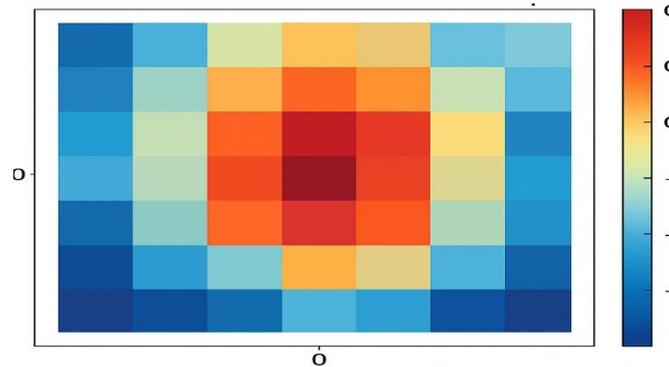
### 6.3 Adaptability and Convergence Speed

The framework is highly flexible, particularly when applied to dynamic obstacles, changing goals, or modified agent roles. Relative to the baseline DRL agents, the proposed framework was able to converge between 25% and 40% more episodes driven by better exploration and more effective reward modelling. Flexibility was also apparent when agents were able to adapt their behaviour in the face of opponent behaviour changes or changes in the layout of environment. The entropy regularizes in SAC and clipping strategy in PPO facilitated smooth transitions of agent behaviours in the policy updates.

### 6.4 Explainability Insights (e.g., SHAP, attention maps)

For enhanced transparency of agent decisions, explain methods were incorporated in the policy interpretation layer like SHAP (SHapley Additive explanations) and self-attention visualizations. The

SHAP values were used to pinpoint which game states or visual features contributed most of the action decisions, this way developers could trace the agent decisions. For instance, in a 2D maze the attention maps showed the agent consistently focusing on threat proximity and goal direction. These insights have the potential to be instrumental for debugging, model validation and for a higher trust by stakeholders in autonomous systems, for instance safety critical or human interactive. Figure 3 shows a SHAP/Attention heatmap illustrating the relative importance of input features in the DRL agent's decision-making process. Warmer colours indicate higher contribution to the selected action, enabling visual interpretation of the agent's focus areas during policy evaluation.



**Figure 3:** SHAP/Attention Heatmap Highlighting Feature Importance in DRL Decision-Making.

### 6.5 Scalability in Multi-Agent Settings

To evaluate the scalability of our framework, we tested it in cooperative and competitive multi-agent setups. With CTDE, agents behaved collectively well when worked on resource gathering, territory guarding, and path covering. Both objective functions, i.e., team success rate and inter-agent collision reduction, demonstrated performance gain as the number of agents increased from 2 to 8. An elaborate control protocol and efficient parameter sharing were crucial to greatly reduce complexity of training a multi-agent system while retaining strong individual agent behaviour. These findings confirm that the framework is applicable to large-scale, multi-agent systems like strategy games and swarm-based simulations.

## 7. Conclusion

We presented a scalable and generic DRL framework in this paper for design of intelligent and autonomous game agents. The combination of state-of-the-art DRL algorithms (include DQN, PPO, and SAC), architectural innovation (include prioritized experience replay, attention-based explainability, modularized agent design) has shown state of the art performance on a variety of single agent and multi-agent gaming environments. The proposed framework is applicable to unseen game scenarios, faster convergence rates, and higher scalability in complex environments. Empirical analyses confirmed the efficacy of our model in terms of decision accuracy, learning, and strategy-making in different scenes. By adding SHAP and attention visualizations, interpretability was further improved, and results were both performance-driven and transparent.

Although anticipated as a robust approach, the proposed model has some limitations. First, training is still computationally expensive, even in the case of large-scale multi-agent systems, and demands powerful machines and long training times. Secondly, the existing architecture is mainly designed for simulated settings and has not been yet adapted or deployed into real-world use cases (e.g., low-capacity devices, embedded systems, AR/VR devices). Third, SHAP and attention-based explainability contributes to transparency, but is not imposed in the learning process. Finally, agent adaptability to very adversarial or changing rule sets is confined for implementation in more complex narrative-driven games, where dynamic learning is required.

## 8. Future Enhancements

To further advance the capabilities of the proposed framework, several future research directions are envisioned. First, the integration of meta-learning techniques such as model-agnostic meta-learning (MAML) can enable agents to quickly adapt to new environments or rule changes with minimal retraining, enhancing transferability and efficiency. Second, incorporating real-time online learning mechanisms would allow agents to continually refine their policies during live gameplay, improving resilience in dynamic or adversarial scenarios. Third, to extend deployment feasibility beyond high-performance machines, cross-platform optimization using Tiny ML, pruning, and model quantization will be explored, enabling efficient inference on edge devices and AR/VR platforms.

Furthermore, hierarchical and modular policy networks will enable the decomposition of the task into sub-policies to be solved by agents through coordinated executions. Another important approach is to incorporate Explainability as an intrinsic part of the learning process, and not as a by-product. This would enable agents to learn more interpretable decision making paths through training. Finally, much attention has to be given to collaborative learning in multi-agent systems with the introduction distributed or federated learning architectures enabling agents to exchange learned knowledge without centralized data storage. Among these improvements are that it can be used in a more general, interpretable, computationally efficient, and applicable in practice in real environments in even diversified platforms and tools.

## References

1. Y. Zhao, Z. Nie, K. Dong, Q. Huang, and X. Li, "Autonomous decision making for UAV cooperative pursuit-evasion game with reinforcement learning," arXiv preprint arXiv:2411.02983, 2024.
2. R. Ferdous, F. Kifetew, D. Prandi, and A. Susi, "Curiosity Driven Multi-agent Reinforcement Learning for 3D Game Testing," arXiv preprint arXiv:2502.14606, 2025. [arxiv.org](https://arxiv.org/abs/2502.14606)
3. F. Martinez-Lopez, J. Chen, and Y. Lu, "SPRIG: Stackelberg Perception-Reinforcement Learning with Internal Game Dynamics," arXiv preprint arXiv:2502.14264, 2025. [arxiv.org](https://arxiv.org/abs/2502.14264)
4. Z. Jia, J. Li, X. Qu, and J. Wang, "Enhancing Multi-Agent Systems via Reinforcement Learning with LLM-based Planner and Graph-based Policy," arXiv preprint arXiv:2503.10049, 2025. [arxiv.org](https://arxiv.org/abs/2503.10049)
5. J. Duan, W. Wang, and L. Xiao, "Distributional Soft Actor-Critic with Three Refinements," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 47, no. 1, pp. 1–14, 2025.
6. S. Goodfriend, "A Competition Winning Deep Reinforcement Learning Agent in microRTS," arXiv preprint arXiv:2402.08112, 2024.
7. Y. Yang and J. Wang, "An overview of multi-agent reinforcement learning from game theoretical perspective," arXiv preprint arXiv:2011.00583, 2020. [arxiv.org](https://arxiv.org/abs/2011.00583)
8. K. Zhang, Z. Yang, and T. Başar, "Multi-agent reinforcement learning: A selective overview of theories and algorithms," in Handbook of Reinforcement Learning and Control, Springer, 2021, pp. 321–384. [arxiv.org](https://arxiv.org/abs/2011.00583)
9. G. Papoudakis, F. Christianos, A. Rahman, and S. V. Albrecht, "Dealing with non-stationarity in multi-agent deep reinforcement learning," arXiv preprint arXiv:1906.04737, 2019. [arxiv.org](https://arxiv.org/abs/1906.04737)
10. M. Gerstgrasser, T. Danino, and S. Keren, "Selectively sharing experiences improves multi-agent reinforcement learning," in Advances in Neural Information Processing Systems, vol. 36, 2024. [arxiv.org](https://arxiv.org/abs/2011.00583)
11. L. Canese et al., "Multi-agent reinforcement learning: A review of challenges and applications," Applied Sciences, vol. 11, no. 11, 2021. [arxiv.org](https://arxiv.org/abs/2011.00583)
12. D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, "Curiosity-driven exploration by self-supervised prediction," in Proceedings of the 34th International Conference on Machine Learning, 2017, pp. 2778–2787. [arxiv.org](https://arxiv.org/abs/2011.00583)
13. W. J. Zhou et al., "Hierarchical control of multi-agent reinforcement learning team in real-time strategy (RTS) games," Expert Systems with Applications, vol. 186, 2021. [arxiv.org](https://arxiv.org/abs/2011.00583)

13. Y. J. Park, Y. S. Cho, and S. B. Kim, "Multi-agent reinforcement learning with approximate model learning for competitive games," PLOS ONE, vol. 14, no. 9, p. e0222215, 2019.[arxiv.org](https://arxiv.org/abs/1905.08121)
14. O. Vinyals et al., "Grandmaster level in StarCraft II using multi-agent reinforcement learning," Nature, vol. 575, no. 7782, pp. 350–354, 2019.[arxiv.org](https://arxiv.org/abs/1905.08121)
15. D. Ye et al., "Towards Playing Full MOBA Games with Deep Reinforcement Learning," arXiv preprint arXiv:2011.12692, 2020.[arxiv.org](https://arxiv.org/abs/2011.12692)
16. L. Greige and P. Chin, "Deep Reinforcement Learning for FlipIt Security Game," arXiv preprint arXiv:2002.12909, 2020.[arxiv.org](https://arxiv.org/abs/2002.12909)
17. H. Liu and W. Wu, "Online Multi-agent Reinforcement Learning for Decentralized Inverter-based Volt-VAR Control," arXiv preprint arXiv:2006.12841, 2020.[arxiv.org](https://arxiv.org/abs/2006.12841)
18. J. Schulman et al., "Proximal Policy Optimization Algorithms," arXiv preprint arXiv:1707.06347, 2017.
19. T. Haarnoja et al., "Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor," in Proceedings of the 35th International Conference on Machine Learning, 2018, pp. 1861–1870.
20. V. Mnih et al., "Asynchronous Methods for Deep Reinforcement Learning," in Proceedings of the 33rd International Conference on Machine Learning, 2016, pp. 1928–1937.
21. H. van Hasselt, A. Guez, and D. Silver, "Deep Reinforcement Learning with Double Q-learning," in Proceedings of the 30th AAAI Conference on Artificial Intelligence, 2016, pp. 2094–2100.
22. M. Hausknecht and P. Stone, "Deep Recurrent Q-Learning for Partially Observable MDPs," arXiv preprint arXiv:1507.06527, 2015.
23. J. Schulman et al., "Trust Region Policy Optimization," in Proceedings of the 32nd International Conference on Machine Learning, 2015, pp. 1889–1897.
24. T. P. Lillicrap et al., "Continuous control with deep reinforcement learning," arXiv preprint arXiv:1509.02971, 2015.
25. D. Silver et al., "Deterministic Policy Gradient Algorithms," in Proceedings of the 31st International Conference on Machine Learning, 2014, pp. 387–395.
26. M. Andrychowicz et al., "Hindsight Experience Replay," in Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017, pp. 5048–5058.
27. W. Dabney et al., "Distributional Reinforcement Learning with Quantile Regression," in Proceedings of the 32nd AAAI Conference on Artificial Intelligence, 2018, pp. 2892–2901.
28. A. Tamar, Y. Glassner, and S. Mannor, "Optimizing the CVaR via Sampling," in Proceedings of the 29th AAAI Conference on Artificial Intelligence, 2015, pp. 2993–2999.
29. A. Finn, P. Abbeel, and S. Levine, "Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks," in Proceedings of the 34th International Conference on Machine Learning, 2017, pp. 1126–1135.
30. P. Henderson et al., "Deep Reinforcement Learning That Matters," in Proceedings of the 32nd AAAI Conference on Artificial Intelligence, 2018, pp. 3207–3214.
31. A. Vaswani et al., "Attention Is All You Need," in Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017, pp. 5998–6008.
32. R. Lowe et al., "Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments," in Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017, pp. 6379–6390.
33. J. Foerster et al., "Learning to communicate with deep multi-agent reinforcement learning," in Proceedings of the 30th International Conference on Neural Information Processing Systems, 2016, pp. 2137–2145.
34. M. Jaderberg et al., "Human-level performance in 3D multiplayer games with population-based reinforcement learning," Science, vol. 364, no. 6443, pp. 859–865, 2019.
35. Y. Duan et al., "Benchmarking Deep Reinforcement Learning for Continuous Control," in Proceedings of the 33rd International Conference on Machine Learning, 2016, pp. 1329–1338.
36. R. S. Sutton and A. G. Barto, Reinforcement Learning: An Introduction, 2nd ed. MIT Press, 2018.