

Reinforcement Learning for Automated Cyber Defense in Dynamic Attack Environments

¹ Atika Nishat, ² Anas Raheem

¹ University of Gurjat, Pakistan

² Air University, Pakistan

Corresponding Email: atikanishat1@gmail.com

Abstract

The rapid evolution of cyber threats and the complexity of dynamic attack environments have rendered traditional rule-based security systems increasingly ineffective. This research explores the application of reinforcement learning (RL) for the development of intelligent, autonomous cyber defense mechanisms capable of adapting in real time to evolving attack strategies. By modeling cybersecurity as a sequential decision-making problem, RL agents learn optimal defense strategies through trial-and-error interactions within simulated environments. We present a comprehensive framework integrating deep reinforcement learning (DRL) with network intrusion detection and response systems to demonstrate the viability of adaptive, automated defense. The research includes the design of training environments using adversarial models, evaluation of agent performance under various attack scenarios, and comparative analysis against static defense systems. Results reveal significant improvements in detection accuracy, response efficiency, and resilience against novel threats, underscoring the potential of RL as a core enabler of future cyber defense automation.

Keywords: Reinforcement Learning, Cybersecurity, Automated Defense, Deep Q-Networks, Dynamic Threat Environments, Intrusion Response, AI for Cyber Defense

I. Introduction

The complexity and dynamism of modern cyber threats demand adaptive defense strategies that go beyond the capabilities of static rule-based systems. Attackers continually evolve their techniques, leveraging zero-day exploits, polymorphic malware, and multi-stage attacks that make detection and mitigation increasingly challenging. Traditional cybersecurity



frameworks, while still essential, often suffer from delayed response, high false-positive rates, and an inability to generalize across previously unseen threats [1]. As the scale and speed of attacks grow, there is a critical need for intelligent systems capable of making autonomous decisions under uncertainty. Reinforcement learning offers a compelling solution to this problem by enabling systems to learn optimal defense strategies through continuous interaction with a simulated or real environment. In RL, an agent learns by receiving rewards or penalties for its actions, gradually improving its behavior to maximize long-term benefits. This paradigm aligns naturally with cyber defense tasks where the environment is often non-deterministic, adversarial, and partially observable. From adaptive firewall management to intrusion response and deception tactics, RL can dynamically adjust policies based on real-time feedback and threat evolution.

Recent advancements in deep reinforcement learning (DRL) have expanded the applicability of RL to high-dimensional state spaces, enabling models to process complex network traffic patterns and system logs. Techniques such as Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and Actor-Critic models have shown promise in cybersecurity applications. However, the integration of these techniques into practical security systems presents unique challenges, including reward engineering, environment simulation, safe exploration, and scalability [2]. The objective of this study is to investigate how DRL can be employed effectively for automated cyber defense in environments characterized by dynamic, evolving attacks. In this paper, we present a framework that models cyber defense as a Markov Decision Process (MDP) and utilize DRL to train agents in a simulated network environment containing diverse attack vectors.

By simulating red-team (attacker) versus blue-team (defender) scenarios, we evaluate the agent's ability to learn proactive and reactive defense strategies. Our contributions include a novel simulation setup, the implementation of multiple RL algorithms, and a comprehensive evaluation of their performance under variable attack scenarios. We also analyze how reward structure and environmental complexity influence learning and generalization. This research aims to bridge the gap between theoretical RL models and their practical deployment in real-world cyber defense infrastructures. By demonstrating tangible benefits in detection, mitigation, and adaptation, we make a case for embedding RL agents within security orchestration, automation, and response (SOAR) platforms [3]. Ultimately, we envision a



future where cybersecurity is powered by continuously learning agents that evolve alongside the threats they are designed to defend against.

II. Methodology

To investigate reinforcement learning in automated cyber defense, we designed a simulated network environment using the Cyber BattleSim toolkit—a Microsoft open-source environment for training RL agents in cybersecurity tasks. The simulation consists of a multihost network topology with various services, vulnerabilities, and attacker behaviors modeled. We incorporated multiple attack vectors such as privilege escalation, lateral movement, data exfiltration, and denial-of-service, allowing for the dynamic representation of real-world attack environments. This sandboxed framework enables safe and controlled training and evaluation of RL agents. The defender agent's task is to detect, respond to, and mitigate these attacks while minimizing resource usage and false positives. The agent observes network states including host logs, traffic summaries, and system alerts, which are encoded into state vectors. Action space includes decisions such as quarantining a host, deploying patches, rerouting traffic, or initiating deception measures. We framed the problem as an MDP with a sparse and delayed reward function where positive reinforcement is given for successful mitigation and negative reinforcement is applied for missed detections or unnecessary actions.

We employed three RL algorithms for comparative analysis: Deep Q-Network (DQN), Advantage Actor-Critic (A2C), and Proximal Policy Optimization (PPO). Each model was trained over 10,000 episodes using \(\varepsilon\)-greedy exploration and experience replay mechanisms where applicable. Neural network architectures were tailored to encode temporal dependencies and environmental correlations, using LSTM layers for sequential log inputs and convolutional layers for network graph embeddings. Hyperparameter tuning was performed using grid search and Bayesian optimization techniques to identify optimal learning rates, discount factors, and batch sizes. To test generalizability, trained agents were exposed to previously unseen attack patterns and adversarial conditions. Metrics for evaluation included detection accuracy, average episode reward, response latency, and false positive rate.



Additionally, to understand the real-time applicability of our models, we deployed the trained agents in a live testbed using virtual machines running vulnerable services. The RL-based defenders interacted with the system through API hooks and log monitors, demonstrating the feasibility of real-time deployment. Logging and audit tools were used to record the agent's behavior, response times, and effectiveness under operational constraints. This methodology not only validates the utility of reinforcement learning in simulated environments but also highlights the pathways to integrating RL agents into existing enterprise security ecosystems. By capturing feedback from real-time operation, agents can continue to learn and adapt, enabling continual defense in an ever-changing threat landscape.

III. Experimental Setup and Results

Our experiments were conducted using two environments: a controlled simulation built with Cyber BattleSim and a live testbed composed of six interconnected virtual machines running Linux-based services. In the simulated environment, we trained each RL agent using an 80/20 train-test split of predefined attack scenarios [4]. Performance metrics were logged every 100 episodes to track convergence and policy improvement. In the live testbed, red-team activities were simulated using the Metasploit framework to test real-world adaptability. In simulation, the DQN agent achieved the highest detection accuracy (91%) after convergence at approximately 7,800 episodes. The PPO agent demonstrated better stability and faster convergence (5,500 episodes) but slightly lower accuracy (88%). A2C lagged behind in both convergence speed and final performance (84%). When tested on novel attack patterns, PPO maintained relatively high performance (86%) due to its ability to generalize across varying inputs, whereas DQN dropped to 78%, indicating potential overfitting to specific attack types.

Latency analysis revealed that PPO responded to threats in under 120ms on average, making it more suitable for real-time scenarios. DQN showed longer response times (~200ms) due to its reliance on replay buffers, while A2C maintained moderate latency (~150ms). False positive rates were lowest in PPO (6%) and highest in A2C (11%), showcasing the importance of balanced exploration strategies in training. In the live testbed, agents responded to attacks like SSH brute force, web server exploitation, and unauthorized lateral



movement. PPO was most resilient in real-world deployments, maintaining consistent performance across test cases. Agents were integrated with a SOAR platform via RESTful APIs, enabling policy execution in a hybrid environment. The agents adapted to increased attack frequency and diversified threats, showing evidence of online learning capabilities through episodic retraining [5].

We also conducted ablation studies to evaluate the impact of reward structure and partial observability [6]. Dense reward schemes led to faster learning but higher false positives, while sparse reward structures promoted more cautious and effective behavior. When observation space was reduced to partial logs, all agents experienced a 10–15% drop in accuracy, emphasizing the need for high-quality monitoring and observability in real deployments. The experimental results collectively demonstrate that reinforcement learning—especially with policy-gradient-based methods—can provide adaptive, robust, and efficient cyber defense mechanisms [7]. However, training stability, safe exploration, and integration into operational environments remain areas requiring careful engineering and continued research.

IV. Discussion

The successful application of reinforcement learning in cyber defense hinges on several interrelated factors including environment modeling, algorithm selection, reward design, and system integration. Our findings indicate that while value-based methods like DQN offer high detection capabilities, their training stability and sensitivity to hyperparameters can be limiting in dynamic, adversarial settings. On the other hand, policy-gradient methods such as PPO present a more stable learning process and better generalization across unseen threats. One major insight from the research is the importance of simulating realistic attack environments. The fidelity of the training environment directly influences the agent's ability to develop transferable skills [8]. Without diverse and evolving adversarial behaviors in the simulation, agents risk becoming brittle or myopic. Integrating red-team behavior, noise injection, and delayed feedback improved the realism and training efficiency of our RL agents [9].



Another critical aspect is the reward engineering process. Sparse and delayed rewards more accurately mimic real-world cyber defense outcomes but increase the learning curve. Our experiments suggest that reward shaping must strike a balance between learning efficiency and behavior generalization [10]. Techniques such as curriculum learning, imitation learning, or reward relabeling may offer future enhancements in this area. The challenge of partial observability and noisy signals is also a barrier to practical deployment. In many real networks, defenders have limited visibility due to encrypted traffic, missing logs, or obfuscated attacks. To address this, future work could integrate recurrent neural networks or transformer-based architectures that can infer hidden states from sequence data, improving robustness under uncertainty.

In operational contexts, integration with existing security infrastructure is crucial. Our deployment in a live testbed demonstrated the feasibility of embedding RL agents within SOAR systems for real-time decision-making [11]. However, ensuring explainability, safety, and human oversight remains essential, especially in high-stakes environments. Approaches such as inverse reinforcement learning and safe policy learning can help ensure alignment with human expectations and legal constraints. Overall, our research confirms the promise of reinforcement learning as a key enabler for adaptive cyber defense. As threats become more automated and intelligent, defense mechanisms must match this evolution with equally dynamic and learning-capable systems. The next frontier lies in federated RL, multi-agent coordination, and continual lifelong learning frameworks that mimic real-world security teams operating at scale and under pressure [12].

V. Conclusion

In conclusion, this research demonstrates that reinforcement learning, particularly when implemented through advanced policy-gradient techniques like PPO, provides a robust and adaptive approach for automated cyber defense in dynamic attack environments. Through carefully designed simulations and live testbed deployments, we showed that RL agents can learn to detect and respond to complex threats with high accuracy, low latency, and strong generalization to novel attacks. While challenges remain in terms of real-time deployment, environment modeling, and reward tuning, our findings establish a strong foundation for



future integration of RL into cybersecurity infrastructures. As cyber threats continue to evolve in sophistication and speed, reinforcement learning offers a scalable, intelligent, and proactive solution for building resilient defense systems capable of adapting and improving autonomously over time.

REFERENCES:

- [1] I. Naseer, "AWS cloud computing solutions: optimizing implementation for businesses," *Statistics, computing and interdisciplinary research,* vol. 5, no. 2, pp. 121-132, 2023.
- [2] I. Naseer, "Cyber defense for data protection and enhancing cyber security networks for military and government organizations," *MZ Computing Journal*, vol. 1, no. 1, pp. 1-8, 2020.
- [3] I. Naseer, "Implementation of Hybrid Mesh firewall and its future impacts on Enhancement of cyber security," *MZ Computing Journal*, vol. 1, no. 2, 2020.
- [4] I. Naseer, "The role of artificial intelligence in detecting and preventing cyber and phishing attacks," *Eur. J. Eng. Sci. Technol*, vol. 11, pp. 82-86, 2024.
- [5] E. Aghaei, X. Niu, W. Shadid, and E. Al-Shaer, "Securebert: A domain-specific language model for cybersecurity," in *International Conference on Security and Privacy in Communication Systems*, 2022: Springer, pp. 39-56.
- [6] I. Naseer, "The crowdstrike incident: Analysis and unveiling the intricacies of modern cybersecurity breaches," *World Journal of Advanced Engineering Technology and Sciences*, vol. 10, p. 3, 2024.
- [7] T. Arif, B. Jo, and J. H. Park, "A Comprehensive Survey of Privacy-Enhancing and Trust-Centric Cloud-Native Security Techniques Against Cyber Threats," *Sensors*, vol. 25, no. 8, p. 2350, 2025.
- [8] I. Naseer, "System malware detection using machine learning for cybersecurity risk and management," *Journal of Science & Technology*, vol. 3, no. 2, pp. 182-188, 2022.
- [9] P. Pandey and A. Patel, "Integrating Security in Cloud-Native Development: A DevSecOps Approach to Resilient Software Systems," in *Data Governance, DevSecOps, and Advancements in Modern Software*: IGI Global Scientific Publishing, 2025, pp. 169-196.
- [10] M. Bayer, T. Frey, and C. Reuter, "Multi-level fine-tuning, data augmentation, and few-shot learning for specialized cyber threat intelligence," *Computers & Security*, vol. 134, p. 103430, 2023.
- [11] I. Naseer, "Machine learning applications in cyber threat intelligence: a comprehensive review," *The Asian Bulletin of Big Data Management*, vol. 3, no. 2, pp. 190-200, 2023.
- [12] E. N. Crothers, N. Japkowicz, and H. L. Viktor, "Machine-generated text: A comprehensive survey of threat models and detection methods," *IEEE Access,* vol. 11, pp. 70977-71002, 2023.