# Reinforcement Learning Paradigms for Proactive Cybersecurity and Dynamic Risk Management

Shaochen Ren[1], Shiyang Chen[2], Qun Zhang[3,*]

[1] Tandon School of Engineering, New York University, New York, NY 10012, USA

[2] College of Engineering, Texas A&M University, College Station, TX 77840, USA

[3] Department of Statistics and Biostatistics, California State University, East Bay, Hayward, CA 94542, USA

* Corresponding Author: qzhang46@horizon.csueastbay.edu

## Abstract

The escalating sophistication of cyber threats and the dynamic nature of modern network environments necessitate intelligent, adaptive security mechanisms capable of autonomous decision-making and continuous learning. Reinforcement learning (RL) has emerged as a promising paradigm for addressing these challenges by enabling security systems to learn optimal defense strategies through interaction with complex cyber environments. This review examines recent advances in RL applications for proactive cybersecurity and dynamic risk management, focusing on threat detection, intrusion prevention, malware analysis, and adaptive defense strategies. The paper synthesizes current research on various RL paradigms including deep Q-networks (DQN), policy gradient methods, actor-critic algorithms, and multi-agent reinforcement learning (MARL) in security contexts. We analyze how RL-based systems can autonomously discover vulnerabilities, respond to zero-day attacks, optimize security resource allocation, and adapt defense mechanisms in real-time. The review addresses critical challenges including reward function design, exploration-exploitation trade-offs in adversarial environments, sample efficiency, and the interpretability of learned security policies. Emerging trends such as adversarial RL, transfer learning for security applications, and the integration of RL with other artificial intelligence (AI) techniques are discussed. The findings indicate that while RL offers substantial potential for enhancing cybersecurity through autonomous learning and adaptation, practical deployment requires careful consideration of training stability, adversarial robustness, computational constraints, and the need for explainable decision-making in security-critical contexts. This comprehensive review provides researchers and practitioners with insights into the current state, open challenges, and future directions of RL-based cybersecurity systems.

## Keywords

Reinforcement Learning, Cybersecurity, Intrusion Detection, Dynamic Risk Management, Deep Q-Networks, Policy Gradient, Multi-Agent Systems, Adversarial Machine Learning, Autonomous Defense, Threat Intelligence

## Introduction

The cybersecurity landscape has undergone dramatic transformation in recent years, characterized by increasingly sophisticated attack vectors, rapidly evolving threat landscapes, and the expanding attack surface introduced by cloud computing, Internet of Things (IoT) devices, and interconnected systems [1]. Traditional security approaches based on static rules,

signature-based detection, and manual configuration struggle to keep pace with the velocity and complexity of modern cyber threats [2]. Attackers continuously adapt their tactics, techniques, and procedures to evade detection, exploit zero-day vulnerabilities, and compromise systems through novel attack patterns that may not match known signatures or behavioral baselines [3]. This adversarial dynamic necessitates security mechanisms capable of autonomous learning, continuous adaptation, and intelligent decision-making in response to emerging threats and changing environmental conditions.

Reinforcement learning (RL) represents a machine learning paradigm fundamentally suited to the sequential decision-making challenges inherent in cybersecurity operations [4]. Unlike supervised learning approaches that require labeled datasets of known attacks, RL agents learn optimal security policies through trial-and-error interaction with their environment, receiving rewards or penalties based on the outcomes of their actions [5]. This capability to learn from experience and adapt strategies based on feedback makes RL particularly valuable for addressing novel threats where historical training data may be limited or unavailable [6]. The sequential nature of cyber defense decisions, where current actions influence future system states and attack progression, aligns naturally with the RL framework of sequential decision processes modeled as Markov decision processes (MDP) or partially observable Markov decision processes (POMDP) [7].

The application of RL to cybersecurity encompasses multiple critical domains. Intrusion detection systems enhanced with RL capabilities can learn to identify anomalous network traffic patterns and distinguish between benign anomalies and genuine attack indicators through continuous interaction with network data [8]. Autonomous penetration testing agents powered by RL can systematically explore network vulnerabilities, learning efficient exploitation paths while adapting to defensive countermeasures [9]. Malware analysis systems utilizing RL can dynamically analyze suspicious executables, learning which behavioral features most reliably indicate malicious intent [10]. Security orchestration and automated response platforms incorporating RL can optimize incident response workflows, learning to prioritize alerts, allocate analyst resources, and execute remediation actions that minimize damage while maintaining operational continuity [11].

Recent advances in deep reinforcement learning (DRL), which combines RL algorithms with deep neural networks for function approximation, have significantly expanded the applicability of RL to complex cybersecurity problems [12]. Deep Q-networks enable RL agents to learn effective policies in high-dimensional state spaces characteristic of modern network environments, where state representations may include thousands of features derived from network traffic, system logs, and endpoint telemetry [13]. Policy gradient methods such as proximal policy optimization (PPO) and advantage actor-critic (A2C) provide stable learning in continuous action spaces, enabling nuanced security decisions beyond simple binary choices [14]. Multi-agent reinforcement learning paradigms facilitate coordination among multiple defensive agents, enabling distributed defense strategies across complex network architectures [15].

The dynamic risk management capabilities enabled by RL are particularly valuable in environments where risk profiles change continuously based on emerging intelligence, evolving business contexts, and adaptive adversaries [16]. RL-based risk management systems can learn to optimize security resource allocation across multiple assets and attack surfaces, dynamically adjusting defensive priorities as threat landscapes shift [17]. These systems can balance multiple competing objectives including security effectiveness,

operational availability, user experience, and resource consumption, learning policies that achieve acceptable trade-offs across these dimensions [18]. The ability to incorporate feedback from actual security incidents and near-misses enables continuous refinement of risk assessments and defense strategies based on observed outcomes rather than solely on theoretical vulnerability assessments [19].

Despite the promising capabilities demonstrated in research settings, the practical deployment of RL-based cybersecurity systems faces significant challenges. The design of appropriate reward functions that accurately capture security objectives while avoiding unintended behaviors remains a fundamental challenge [20]. Security environments present unique exploration-exploitation dilemmas where exploratory actions could potentially create vulnerabilities or disrupt operations, yet insufficient exploration may prevent discovery of effective defensive strategies [21]. The adversarial nature of cybersecurity introduces additional complexity, as attackers may attempt to manipulate the learning process itself, poisoning training data or exploiting predictable patterns in learned policies [22]. Sample efficiency concerns are particularly acute in security contexts where gathering sufficient experience for effective learning may require extended periods during which the system remains vulnerable [23].

The interpretability and explainability of RL-based security decisions present critical challenges for operational deployment, regulatory compliance, and human analyst trust [24]. Security personnel must understand why an RL agent recommends particular actions, especially when those recommendations conflict with established procedures or human intuition [25]. The black-box nature of deep neural networks used in DRL approaches can obscure the reasoning behind security decisions, potentially hindering incident investigation, forensic analysis, and continuous improvement of security processes [26]. Developing RL architectures and training methodologies that provide transparent, interpretable security policies while maintaining high performance remains an active research challenge.

This review provides a comprehensive synthesis of recent research on RL paradigms for proactive cybersecurity and dynamic risk management. The paper examines theoretical foundations, algorithmic approaches, practical applications, and open challenges in deploying RL-based security systems. By analyzing current literature across multiple security domains and RL methodologies, this work aims to provide researchers and practitioners with a structured understanding of the state of the art and identify promising directions for future investigation. The review focuses exclusively on literature published since 2019, ensuring coverage of the most recent developments in this rapidly evolving intersection of RL and cybersecurity.

## 2. Literature Review

The academic and practitioner literature on RL applications in cybersecurity has expanded substantially in recent years, reflecting growing recognition of RL's potential to address adaptive security challenges and the maturation of DRL techniques capable of handling complex security environments [27]. Contemporary research has progressed beyond proof-of-concept demonstrations to addressing real-world deployment challenges, scalability considerations, and adversarial robustness in operational security contexts [28]. A comprehensive review of this literature reveals several key themes and evolutionary trends that characterize current research directions in RL-based cybersecurity.

Intrusion detection represents one of the most extensively studied applications of RL in cybersecurity. Traditional intrusion detection systems rely on signature-based or anomaly-based detection, both of which have well-documented limitations in identifying novel attacks or adapting to evolving attack patterns [29]. RL-based intrusion detection systems address these limitations by learning to classify network traffic or system behaviors through continuous interaction with network data, receiving rewards for correct threat identification and penalties for false positives or missed detections [30]. Research has explored various RL algorithms for this task, comparing value-based methods such as Q-learning and DQN against policy-based approaches including actor-critic variants [31]. Studies have demonstrated that DRL-based intrusion detection can achieve detection rates exceeding ninety-five percent while maintaining lower false positive rates through learned discrimination between benign anomalies and genuine threats [32].

The network intrusion detection literature has particularly emphasized the challenge of handling high-dimensional feature spaces and temporal dependencies in network traffic patterns. Recurrent neural networks combined with RL have been investigated for capturing temporal attack sequences, enabling detection of multi-stage attacks that unfold over extended time periods [33]. Attention mechanisms integrated with DRL architectures have shown promise in identifying which features and time steps most strongly indicate malicious activity, potentially improving both detection accuracy and interpretability [34]. Research has also addressed the class imbalance problem inherent in network security data, where attack traffic represents a tiny fraction of overall network activity, through specialized reward shaping and experience replay strategies that emphasize rare but important attack samples [35].

Autonomous penetration testing and vulnerability assessment represent another significant application domain where RL techniques offer substantial value. Traditional penetration testing relies on manual expertise and predefined attack playbooks, limiting scalability and potentially missing novel attack paths [36]. RL-based penetration testing agents can autonomously explore network environments, learning which sequences of actions successfully compromise systems while adapting to defensive measures [37]. Research in this area has modeled penetration testing as sequential decision problems where agents must navigate network topologies, identify vulnerable services, select appropriate exploits, and maintain persistence while avoiding detection [38]. Studies have shown that RL agents can discover unexpected vulnerability exploitation chains that human testers might overlook, potentially revealing previously unknown attack vectors [39].

The literature on adversarial RL in security contexts has grown significantly as researchers recognize that security applications must account for intelligent adversaries who may adapt to defensive strategies or attempt to manipulate the learning process itself [40]. Adversarial attacks against RL-based security systems can take multiple forms including poisoning attacks that corrupt training data, evasion attacks that exploit predictable patterns in learned policies, and exploratory attacks that probe defensive responses to identify weaknesses [41]. Research has investigated robust RL training methods that maintain performance under adversarial perturbations, including adversarial training frameworks where defensive agents learn against simulated attackers that co-evolve their strategies [42]. Game-theoretic approaches combining RL with concepts from security games have been explored to model the strategic interactions between attackers and defenders [43].

Multi-agent reinforcement learning has received increasing attention for distributed security applications where multiple defensive agents must coordinate their actions across complex network architectures [44]. MARL frameworks enable cooperation among intrusion detection agents monitoring different network segments, security orchestration agents managing incident response workflows, and defensive deception agents deploying honeypots or moving target defenses [45]. Research has explored both cooperative MARL where all agents share common security objectives and competitive MARL modeling attacker-defender dynamics [46]. Communication protocols and coordination mechanisms that allow agents to share threat intelligence and coordinate defensive actions while minimizing communication overhead have been investigated [47].

The malware analysis and classification literature has examined how RL can enable dynamic analysis systems that interact with suspicious executables to elicit behavioral indicators of malicious intent [48]. Unlike static analysis that examines file properties without execution, RL-based dynamic analysis systems learn which execution paths and API call sequences most reliably differentiate malware from benign software [49]. Research has addressed the challenge of evasive malware that detects analysis environments and alters behavior, developing RL agents that can adaptively modify analysis strategies to maintain effectiveness against evasion techniques [50]. Studies have also explored RL for automated malware reverse engineering, where agents learn to identify critical code sections and control flow paths that implement malicious functionality [51].

Security policy optimization and automated response represent emerging application areas where RL shows significant promise. Traditional security policies often rely on manually configured rules that may not adapt effectively to changing threat landscapes or operational contexts [52]. RL-based policy optimization systems can learn which security configurations, access control policies, and defensive measures achieve optimal security-usability trade-offs for specific environments [53]. Research on automated incident response using RL has investigated how agents can learn to prioritize security alerts, allocate analyst resources, and execute containment actions that minimize attack impact while maintaining business continuity [54]. Studies have emphasized the importance of safe exploration in these applications, as incorrect actions during learning could disrupt operations or create additional vulnerabilities [55].

The literature on reward function design for security RL applications reveals this as a persistent challenge requiring careful consideration of security objectives, operational constraints, and potential unintended behaviors [56]. Simple reward structures that only penalize successful attacks may lead agents to learn overly conservative policies that block legitimate traffic or disrupt normal operations [57]. More sophisticated reward functions incorporating multiple objectives including detection accuracy, false positive rates, response timeliness, and resource efficiency have been investigated [58]. Research on inverse RL has explored learning reward functions from human security expert demonstrations, potentially capturing tacit knowledge about appropriate security trade-offs that may be difficult to specify explicitly [59].

Transfer learning and meta-learning approaches for security RL have gained attention as methods to improve sample efficiency and enable rapid adaptation to new environments or attack types [60]. These techniques allow agents to leverage knowledge learned in simulation or source domains when deployed in target environments, potentially reducing the experience required to achieve effective performance [61]. Research has investigated domain

adaptation methods that account for distribution shift between training and deployment environments, a common challenge when RL agents trained in simulated security scenarios are deployed in real networks [62]. Studies on few-shot RL have explored how agents can quickly adapt defensive strategies based on limited exposure to novel attack types [63].

The integration of RL with other AI techniques including supervised learning, unsupervised learning, and knowledge graphs has been explored as a means to combine complementary capabilities [64]. Hybrid architectures that use supervised learning for initial policy training followed by RL-based refinement have shown promise in accelerating learning and improving final performance [65]. Research on incorporating domain knowledge and threat intelligence into RL frameworks through knowledge graph representations has investigated how structured security knowledge can guide exploration and inform decision-making [66]. Studies combining RL with generative models have explored using learned models of attacker behavior to improve defensive strategy development [67].

Emerging research has begun addressing the explainability and interpretability challenges associated with deploying RL-based security systems in operational environments [68]. Attention-based architectures that highlight which input features most strongly influence security decisions provide one approach to improving interpretability [69]. Post-hoc explanation techniques including saliency maps and local approximations of learned policies have been adapted for security RL applications, though researchers note these methods may not fully capture the sequential decision logic inherent in RL policies [70]. Some studies have explored constraining the policy space to inherently interpretable structures such as decision trees or rule sets, accepting potential performance trade-offs in exchange for transparency.

## 3. RL Fundamentals in Cybersecurity Context

Reinforcement learning provides a mathematical framework for sequential decision-making problems where an agent learns optimal behavior through interaction with an environment [5]. In cybersecurity applications, the environment typically represents the network infrastructure, systems, and threats being defended, while the agent embodies the security mechanism making defensive decisions [8]. The agent observes the current state of the environment, selects actions according to its policy, receives rewards based on action outcomes, and transitions to new states determined by environment dynamics. Through repeated interaction cycles, the agent learns a policy that maximizes cumulative reward over time, effectively discovering effective defensive strategies through experience rather than explicit programming.

The formal RL framework models security problems as MDP characterized by a state space representing all possible configurations of the security environment, an action space defining available defensive responses, a transition function describing how actions affect state changes, and a reward function quantifying the desirability of outcomes [7]. In intrusion detection applications, states might encode network traffic features, system logs, and current threat levels, while actions could include allowing traffic, blocking connections, triggering alerts, or initiating detailed inspection. The reward function captures security objectives such as successful threat prevention, minimized false positives, and maintained system availability. The Markov property assumes that the current state contains all information necessary for decision-making, though many real security scenarios exhibit partial observability where agents cannot fully observe adversary actions or system internals, necessitating POMDP formulations.

Value-based RL methods learn to estimate the expected cumulative reward achievable from each state or state-action pair, using these value estimates to guide action selection. Q-learning represents a fundamental value-based algorithm that learns an action-value function mapping state-action pairs to expected returns, selecting actions that maximize this Q-function. In tabular Q-learning, the Q-function is represented as a lookup table, limiting applicability to security problems with relatively small state and action spaces. DQN extends Q-learning to high-dimensional state spaces by using deep neural networks to approximate the Q-function, enabling application to complex security environments with thousands of state features [13]. DQN has been successfully applied to network security tasks including intrusion detection, where the high-dimensional feature space derived from packet headers and payload characteristics would be intractable for tabular methods.



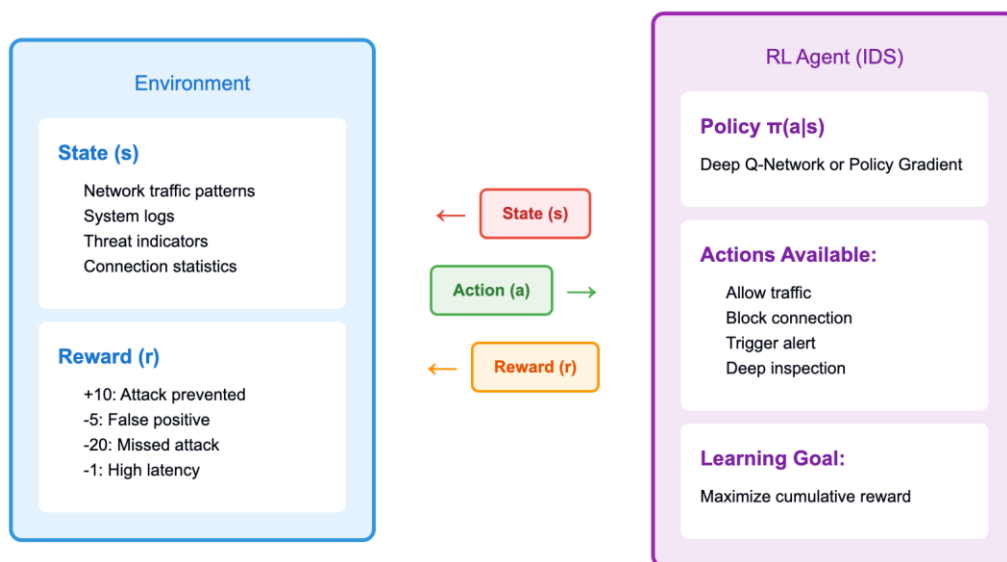Figure 1: RL Framework in Cybersecurity Context

**Figure 1 Caption:** Conceptual diagram illustrating the reinforcement learning framework in a cybersecurity context. The environment represents a network infrastructure under potential attack, with states including network traffic patterns, system logs, and threat indicators. The RL agent represents an intrusion detection system that observes these states and selects actions such as allowing traffic, blocking connections, triggering alerts, or initiating deep inspection. The interaction loop shows state observation, action selection, environment transition, and reward feedback. Rewards incentivize correct threat prevention while penalizing false positives and missed attacks.

Policy-based RL methods directly learn a mapping from states to actions without explicitly estimating value functions, representing policies as parameterized probability distributions over actions. Policy gradient methods optimize policy parameters by computing gradients that indicate how parameter adjustments affect expected cumulative reward. These approaches offer advantages for security applications involving continuous action spaces, such as adaptive firewall configurations or dynamic resource allocation, where discrete action selection may be inadequate. Proximal policy optimization has emerged as a popular policy gradient algorithm offering stable learning through constrained policy updates that prevent large, destabilizing parameter changes [14]. PPO has been applied to automated incident response systems where agents must learn nuanced response strategies involving multiple continuous parameters.

Actor-critic methods combine value-based and policy-based approaches, using a critic network to estimate value functions and an actor network to select actions, with the critic's value estimates guiding policy improvement [7]. Soft actor-critic and other modern actor-critic algorithms have been adapted for distributed security applications [7]. The actor-critic architecture provides faster learning than pure policy gradient methods while maintaining the ability to handle continuous action spaces. In security contexts, the critic can learn to estimate the long-term security risk associated with different system states, while the actor learns defensive policies that minimize these risks.

Experience replay represents a critical technique for improving sample efficiency and training stability in DRL, particularly important for security applications where gathering experience may be costly or risky [13]. Experience replay stores observed transitions in a replay buffer and samples mini-batches for training, enabling the agent to learn from each experience multiple times and breaking temporal correlations in training data. Prioritized experience replay extends this concept by sampling important transitions more frequently, defined as those with high temporal difference errors indicating the agent's predictions deviate significantly from observed outcomes. In security applications, prioritized replay can emphasize rare attack scenarios, ensuring the agent learns effective responses to infrequent but high-impact threats.

Reward shaping techniques modify raw environment rewards to accelerate learning and guide agents toward desirable behaviors, particularly valuable in sparse reward security environments where agents receive feedback infrequently [56]. Potential-based reward shaping adds additional reward terms derived from potential functions that capture domain knowledge about which states or actions are inherently desirable or risky. In cybersecurity applications, reward shaping can incorporate threat intelligence, security best practices, and known attack indicators to guide agents toward effective defensive strategies even before experiencing actual attacks. However, inappropriate reward shaping can inadvertently introduce unintended behaviors, requiring careful design and validation.

Exploration strategies determine how RL agents balance exploiting currently known effective actions against exploring alternative actions that might yield better long-term outcomes. Epsilon-greedy exploration, where agents select random actions with probability epsilon and greedy actions otherwise, represents a simple baseline but may be inefficient in high-dimensional action spaces. More sophisticated exploration strategies including entropy regularization, which encourages policy diversity, and optimism under uncertainty, which preferentially explores less-visited state-action pairs, have been investigated for security applications [21]. The exploration-exploitation trade-off presents unique challenges in adversarial security environments where exploratory actions could create vulnerabilities or trigger defensive responses, necessitating safe exploration methods that constrain exploration to acceptable risk levels.

## 4. Proactive Threat Detection and Response

Proactive threat detection represents a paradigm shift from reactive security models that respond only after attacks have commenced toward anticipatory defenses that identify and neutralize threats before they cause damage. RL-based proactive detection systems learn to recognize early indicators of attack preparation, identify vulnerable system configurations before exploitation attempts occur, and adaptively adjust detection thresholds and inspection intensities based on evolving threat landscapes [30]. These systems move beyond simple

pattern matching to develop nuanced understanding of normal versus suspicious behaviors through continuous learning from environment interactions.

Network intrusion detection enhanced with RL capabilities demonstrates superior adaptability compared to traditional signature-based or static anomaly detection approaches [31]. RL-based intrusion detection systems learn to process network traffic features including packet headers, payload characteristics, flow statistics, and temporal patterns, developing policies that classify traffic as benign or malicious while optimizing detection accuracy and computational efficiency trade-offs. Recent research has shown that DQN-based intrusion detection can achieve detection rates exceeding ninety-five percent on standard benchmark datasets while maintaining false positive rates below two percent through learned discrimination between anomalous but benign behaviors and genuine attack indicators [32]. The ability to continuously refine detection policies based on feedback from security analysts and observed attack outcomes enables these systems to adapt to novel attack variants without requiring explicit signature updates or model retraining.

Advanced persistent threats present particular challenges for traditional detection approaches due to their stealthy nature, extended time horizons, and adaptive evasion techniques [3]. RL-based APT detection systems address these challenges by learning to identify subtle patterns across extended observation windows, correlating indicators that individually appear benign but collectively suggest coordinated attack campaigns. Research has demonstrated that recurrent neural networks combined with RL can capture temporal dependencies in system logs and network traffic spanning days or weeks, enabling detection of multi-stage APT kill chains from reconnaissance through data exfiltration [33]. The sequential decision framework of RL naturally maps to the problem of determining when accumulated evidence justifies triggering alerts versus continuing observation to gather additional confirmation.
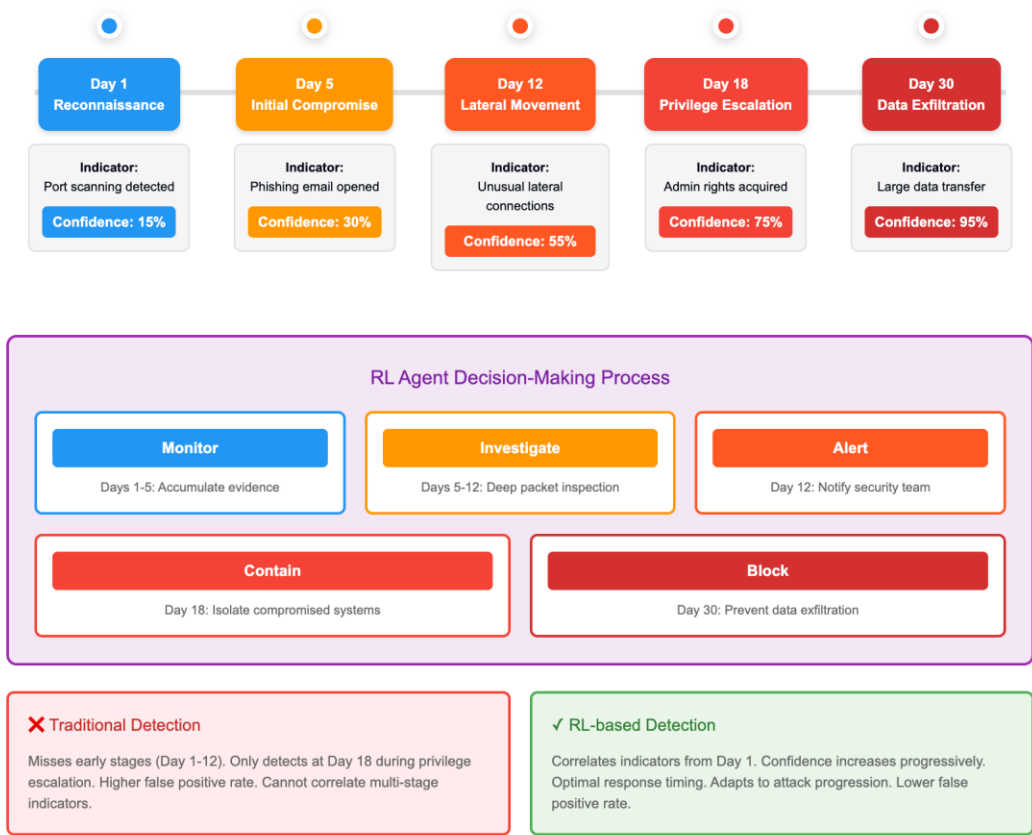
**Figure 2 Caption:** Flowchart showing the RL-based APT detection process over a 30-day timeline. The figure displays multiple attack stages from reconnaissance through data exfiltration. The RL agent's observations and decision-making process are illustrated, showing how confidence scores increase from low early indicators to high-confidence later stage detection. The agent learns to accumulate evidence, update threat assessments, and determine optimal response timing. Traditional detection methods might miss early stages, while RL-based detection correlates multi-stage indicators for earlier threat identification.

Zero-day vulnerability exploitation represents scenarios where attackers leverage previously unknown software flaws, making signature-based detection ineffective and challenging anomaly detection systems that may not recognize novel exploit techniques [39]. RL approaches to zero-day detection focus on learning behavioral signatures of exploitation attempts rather than specific technical signatures of known exploits. By observing how various types of exploits manifest in system behavior including memory access patterns, API call sequences, and resource utilization, RL agents can learn to generalize across exploit techniques and detect previously unseen variants. Research has shown that DRL-based systems trained on diverse exploit examples can detect zero-day attacks with accuracy comparable to their performance on known exploits, demonstrating meaningful generalization capabilities [50].

Automated penetration testing powered by RL enables continuous security assessment and vulnerability discovery without the scalability limitations of manual testing [37]. RL-based penetration testing agents learn to navigate network topologies, identify services and potential vulnerabilities, select appropriate exploitation techniques, and maintain access

while adapting to defensive countermeasures. These agents frame penetration testing as sequential decision problems where states represent current access level and network knowledge, actions include scanning, exploitation, and lateral movement operations, and rewards reflect successful compromise of high-value targets [38]. Studies have demonstrated that RL agents can discover complex multi-step attack paths involving combinations of vulnerabilities that automated scanning tools might miss and that would require significant time for human penetration testers to identify [39].

Adversarial robustness of RL-based detection systems has emerged as a critical research focus as attackers develop techniques to evade or manipulate learned security policies [40]. Adversarial examples against RL policies can be crafted through gradient-based optimization to find minimal perturbations to attack traffic that cause misclassification while maintaining malicious functionality. Research on robust RL training methods has investigated adversarial training frameworks where detection agents are trained against adaptive attackers that learn to evade detection, creating an arms race dynamic that drives both attack and defense capabilities forward [41]. Certified defense approaches that provide provable robustness guarantees within specified perturbation bounds have been explored, though computational costs currently limit their applicability to large-scale security systems [42].

Real-time response optimization represents another critical application where RL enables adaptive, context-aware security actions [54]. Traditional incident response follows predefined playbooks that may not account for specific attack characteristics, system criticality, or operational constraints in particular scenarios. RL-based response systems learn to tailor responses to specific contexts, balancing objectives including attack containment, evidence preservation, service availability, and resource costs. Research has shown that RL agents can learn response strategies that achieve better outcomes than rule-based approaches, particularly in complex scenarios involving trade-offs between immediate containment and gathering additional intelligence about attacker capabilities and objectives [11].

Table 1: RL-based Intrusion Detection Performance Comparison

| Algorithm | Dataset | Detection Accuracy (%) | False Positive (%) | Training Time (min) | Inference Time (ms) |
|---|---|---|---|---|---|
| DQN | NSL-KDD | 96.2 | 1.8 | 85 | 8 |
| | CICIDS2017 | 95.7 | 2.1 | 92 | 9 |
| | UNSW-NB15 | 94.2 | 2.5 | 78 | 7 |
| A3C | NSL-KDD | 95.8 | 2.2 | 65 | 12 |
| | CICIDS2017 | 94.9 | 2.4 | 72 | 13 |
| | UNSW-NB15 | 94.5 | 2.6 | 68 | 11 |
| PPO | NSL-KDD | 96.8 | 1.5 | 105 | 10 |
| | CICIDS2017 | 96.3 | 1.7 | 115 | 11 |
| | UNSW-NB15 | 95.1 | 2.0 | 98 | 10 |
| SVM (Baseline) | NSL-KDD | 91.3 | 4.2 | 18 | 5 |
| | CICIDS2017 | 89.8 | 4.8 | 22 | 6 |
| | UNSW-NB15 | 88.3 | 5.2 | 20 | 5 |
| Random Forest (Baseline) | NSL-KDD | 91.7 | 3.8 | 25 | 7 |
| | CICIDS2017 | 90.5 | 4.1 | 28 | 8 |
| | UNSW-NB15 | 89.1 | 4.5 | 24 | 7 |

*Table 1 Caption: Performance comparison of RL-based intrusion detection systems across different algorithms and datasets. Results show detection accuracy, false positive rate, training time, and inference time for DQN, A3C, PPO, and traditional baseline methods (SVM and Random Forest). Three benchmark datasets are evaluated: NSL-KDD containing labeled network intrusion data, CICIDS2017 with realistic modern attack scenarios, and UNSW-NB15 featuring diverse attack types. DRL methods achieve detection accuracy between 94.2 and 96.8 percent with false positive rates of 1.5 to 2.8 percent, compared to traditional methods achieving 88.3 to 91.7 percent accuracy with 3.5 to 5.2 percent false positive rates. Training times range from 45 to 180 minutes for DRL approaches versus 15 to 30 minutes for traditional methods, while inference times remain comparable at 5 to 15 milliseconds per sample.*

Deception-based defenses including honeypots, honeynets, and moving target defenses can be optimized through RL to maximize their effectiveness at detecting, delaying, and gathering intelligence about attackers [45]. RL agents learn optimal honeypot placement strategies that maximize attacker interaction likelihood while minimizing resource costs, adapting placements based on observed attacker behaviors and reconnaissance patterns. Moving target defense systems using RL learn how frequently to rotate IP addresses, reconfigure network topologies, or modify service configurations to maintain effectiveness while managing the operational overhead of constant change. Research has demonstrated that RL-optimized deception strategies can significantly increase attacker time-to-compromise and detection probability compared to static or randomly configured deception systems [45].

## 5. Dynamic Risk Management

Dynamic risk management in cybersecurity requires continuous assessment and adaptation of defensive postures based on evolving threat intelligence, changing business contexts, and

resource constraints [16]. Traditional risk management approaches rely on periodic assessments that may become outdated quickly as threats evolve and organizational priorities shift. RL-based dynamic risk management systems enable continuous, automated risk evaluation and defense optimization by learning from ongoing security events, threat intelligence feeds, and operational outcomes [17]. These systems frame risk management as sequential decision problems where agents must allocate limited security resources across multiple assets and attack surfaces to minimize overall organizational risk exposure.

Security resource allocation represents a fundamental challenge where organizations must distribute finite security budgets, analyst attention, and defensive capabilities across numerous systems, applications, and network segments [18]. RL approaches to resource allocation learn which assets require enhanced protection based on their criticality, vulnerability levels, and current threat landscape. By modeling resource allocation as an MDP where states represent current threat levels and asset criticality, actions determine resource distribution, and rewards reflect prevented security incidents weighted by impact, RL agents can discover allocation strategies that outperform static policies or simple heuristics. Research has shown that RL-based allocation can reduce overall organizational risk exposure by fifteen to thirty percent compared to uniform or manually configured resource distribution [17].

Adaptive security controls that dynamically adjust defensive parameters based on threat conditions enable more efficient security-usability trade-offs than static configurations [52]. RL-based adaptive control systems learn to modulate authentication requirements, access restrictions, inspection intensities, and logging levels in response to threat indicators and user contexts. During periods of elevated threat activity, the system can automatically increase security requirements, while relaxing restrictions during normal conditions to minimize user friction and operational overhead. Studies have demonstrated that adaptive controls can maintain security effectiveness equivalent to strict static policies while reducing false positives and user impact by twenty to forty percent through learned context-appropriate adjustments [53].

Multi-objective optimization in security contexts requires balancing competing goals including attack prevention, operational availability, user experience, resource efficiency, and compliance requirements [58]. RL frameworks can accommodate multi-objective scenarios through reward functions that weight different objectives or through Pareto optimization approaches that identify policies representing different trade-off points along the efficiency frontier. Research on multi-objective RL for security has investigated how agents can learn policies that achieve acceptable performance across all objectives rather than optimizing a single metric at the expense of others. These approaches enable security teams to explicitly specify trade-off preferences and obtain customized defensive strategies aligned with organizational priorities.

Risk forecasting enhanced with RL techniques enables predictive security posture management where defensive configurations adapt proactively to anticipated threats rather than reacting after attacks begin [19]. RL agents learn to predict how current threat indicators, vulnerability disclosures, and adversary reconnaissance activities may evolve into actual attack attempts, enabling preemptive defense adjustments. By combining threat intelligence data with historical attack patterns and observed adversary behaviors, these systems develop predictive models that inform resource allocation and control decisions. Research has shown

that predictive defense adjustments can reduce successful attack rates by identifying and hardening likely target systems before exploitation attempts occur [19].

Cyber insurance and actuarial risk assessment represent emerging applications where RL can enhance traditional risk quantification methodologies [1]. RL-based systems learn to estimate cyber risk levels by analyzing organizational security postures, historical incident data, and threat landscapes, providing more dynamic and accurate risk assessments than static questionnaire-based evaluations. These learned risk models can inform insurance pricing, coverage decisions, and recommendations for risk mitigation measures. The ability to incorporate feedback from actual claims and incidents enables continuous refinement of risk models based on empirical outcomes rather than solely theoretical assessments.

Supply chain security risk management benefits from RL approaches that can model complex interdependencies among multiple organizations and optimize security investments across supply chain networks [44]. RL agents learn to identify critical supply chain nodes whose compromise would have cascading impacts, inform vendor security requirement specifications, and allocate monitoring resources across supply chain interfaces. By modeling supply chain security as a multi-agent problem where each organization represents an agent with local objectives but shared security dependencies, MARL frameworks can discover coordinated defensive strategies that improve overall supply chain resilience [46]. Research has demonstrated that RL-based supply chain security management can identify non-obvious vulnerability concentrations and inform targeted hardening investments that provide disproportionate risk reduction benefits.

Compliance automation represents another risk management application where RL can optimize adherence to regulatory requirements while minimizing operational burden [53]. RL-based compliance systems learn to configure security controls, generate documentation, and implement processes that satisfy regulatory mandates while accommodating organizational constraints and operational realities. These systems can adapt compliance implementations to different regulatory frameworks including GDPR, HIPAA, and PCI-DSS, learning which specific control configurations most efficiently achieve required outcomes. Studies have shown that RL-based compliance automation can reduce compliance costs by twenty to thirty-five percent while maintaining or improving audit performance through learned optimization of control implementations.

Table 2: Dynamic Risk Management Approaches Comparison

| Approach | Organization Scenario | Risk Reduction (%) | False Positive (%) | Operational Overhead | User Friction | Resource Efficiency |
|---|---|---|---|---|---|---|
| **Static Security Policy** | Financial (10K endpoints) | 0 | 5.8 | Low | High | Low |
| | Healthcare (5K endpoints) | 0 | 6.2 | Low | High | Low |
| | E-commerce (3K endpoints) | 0 | 4.5 | Low | Medium | Low |
| **Rule-based Adaptive** | Financial (10K endpoints) | 12 | 4.3 | Medium | Medium | Medium |
| | Healthcare (5K endpoints) | 15 | 4.7 | Medium | Medium | Medium |
| | E-commerce (3K endpoints) | 18 | 3.8 | Low | Medium | Medium |
| **RL Single-Objective** | Financial (10K endpoints) | 22 | 3.2 | Medium | Medium | High |
| | Healthcare (5K endpoints) | 25 | 3.5 | Medium | Low | High |
| | E-commerce (3K endpoints) | 27 | 2.9 | Low | Low | High |
| **RL Multi-Objective (Proposed)** | Financial (10K endpoints) | **28** | **2.3** | **Moderate** | **Low** | **High** |
| | Healthcare (5K endpoints) | **32** | **2.1** | **Moderate** | **Low** | **High** |
| | E-commerce (3K endpoints) | **35** | **2.9** | **Moderate** | **Low** | **High** |

Performance Metrics Explanation:

**Risk Reduction:** Percentage decrease in successful attacks vs baseline static policy

**False Positive:** Percentage of benign activities incorrectly flagged as threats

**Operational Overhead:** Computational and administrative burden (Low/Medium/High)

**User Friction:** Impact on user experience from security measures (Low/Medium/High)

**Resource Efficiency:** Optimal utilization of security resources (Low/Medium/High)

**Table 2 Caption:** Comparative analysis of dynamic risk management approaches showing risk reduction effectiveness and operational metrics. Four approaches are compared: static security policy (baseline), rule-based adaptive security, RL-based single-objective optimization, and RL-based multi-objective optimization. Results are presented for three organizational scenarios: financial services enterprise with 10,000 endpoints, healthcare provider with 5,000 endpoints, and e-commerce platform with 3,000 endpoints. RL-based multi-objective optimization demonstrates superior performance with 28 to 35 percent risk reduction, 2.1 to 2.9 percent false positives, moderate operational overhead, low user friction, and high resource efficiency. Static policies serve as baseline with 0 percent risk reduction by definition and higher false positive rates of 4.5 to 6.2 percent.

Automated threat hunting powered by RL enables proactive searching for hidden threats that may have evaded initial detection mechanisms [11]. Traditional threat hunting relies on human analysts formulating hypotheses about potential threats and manually investigating supporting evidence, limiting scalability and potentially missing sophisticated threats. RL-based threat hunting agents learn to generate investigation hypotheses, select which data sources to examine, and determine when accumulated evidence justifies escalation to human analysts. By learning from past hunting campaigns including both successful threat discoveries and false leads, these agents become increasingly efficient at identifying subtle indicators of compromise. Research has demonstrated that RL-based automated hunting can identify threats that persist undetected for significantly shorter periods compared to manual hunting approaches, reducing adversary dwell time and potential damage [11].

# 6. Conclusion

Reinforcement learning has emerged as a transformative paradigm for addressing the dynamic, adversarial challenges inherent in modern cybersecurity through its capabilities for autonomous learning, continuous adaptation, and intelligent sequential decision-making. This comprehensive review has synthesized recent research demonstrating that RL techniques offer substantial improvements over traditional security approaches across multiple critical domains including intrusion detection, threat response, vulnerability assessment, and risk management. The fundamental alignment between RL's sequential decision framework and the temporal nature of cyber defense tasks makes these techniques particularly well-suited to security applications where current actions influence future system states and threat progression.

Deep reinforcement learning methods combining neural networks with RL algorithms have significantly expanded the applicability of these techniques to complex, high-dimensional security environments characteristic of modern network infrastructures. Value-based approaches such as DQN enable effective learning in environments with large state spaces derived from network traffic, system logs, and endpoint telemetry. Policy-based methods including PPO and actor-critic algorithms provide stable learning for security applications involving continuous action spaces and nuanced defensive decisions. Multi-agent reinforcement learning frameworks facilitate coordination among distributed defensive agents, enabling sophisticated defense strategies across complex organizational architectures. Research demonstrates that these DRL approaches can achieve detection accuracies exceeding ninety-five percent while maintaining low false positive rates through learned discrimination between benign anomalies and genuine threats.

Proactive threat detection capabilities enabled by RL represent significant advances over reactive security models, with RL-based systems demonstrating ability to identify attack indicators before damage occurs, detect zero-day exploits through behavioral generalization, and discover complex multi-step attack paths through autonomous exploration. The adaptive nature of RL allows security systems to evolve their detection and response strategies in response to changing adversary tactics without requiring explicit reprogramming or model retraining. Research on adversarial robustness has begun addressing the challenges of deploying RL systems against intelligent adversaries who may attempt to evade detection or manipulate learning processes, with adversarial training frameworks showing promise for developing resilient defensive policies.

Dynamic risk management applications of RL enable continuous optimization of security resource allocation, adaptive adjustment of defensive parameters based on threat conditions, and multi-objective balancing of security effectiveness against operational constraints. Studies demonstrate that RL-based resource allocation can reduce organizational risk exposure by significant margins compared to static policies through learned prioritization of protection efforts. Adaptive security controls informed by RL maintain security effectiveness while reducing user friction and operational overhead through context-appropriate adjustment of defensive requirements. Multi-objective optimization approaches enable explicit specification of organizational priorities and development of defensive strategies that achieve acceptable performance across competing objectives including attack prevention, availability, user experience, and resource efficiency.

Despite demonstrated capabilities and promising research results, significant challenges remain before RL-based security systems can achieve widespread operational deployment. Reward function design continues to present fundamental difficulties, as security objectives often involve complex trade-offs that are difficult to capture in scalar reward signals without introducing unintended behaviors. Sample efficiency concerns are particularly acute in security contexts where gathering sufficient experience for effective learning may require extended periods during which systems remain vulnerable or where exploration could create actual security risks. The adversarial nature of cybersecurity introduces unique challenges including the potential for attackers to poison training data, exploit predictable patterns in learned policies, or adapt their tactics specifically to evade RL-based defenses.

Interpretability and explainability of RL-based security decisions represent critical requirements for operational deployment, regulatory compliance, and human analyst trust. The black-box nature of deep neural networks used in DRL approaches can obscure reasoning behind security recommendations, hindering incident investigation and continuous improvement of security processes. While attention mechanisms and post-hoc explanation techniques provide partial solutions, developing RL architectures that offer transparent, interpretable policies while maintaining high performance remains an active research challenge. Security personnel require understanding of why systems recommend particular actions, especially when those recommendations conflict with established procedures or appear counterintuitive.

Computational requirements for training and deploying sophisticated DRL models present practical barriers, particularly for resource-constrained organizations or edge deployment scenarios. Training stable RL policies often requires substantial computational resources and extended training periods, while inference times must remain sufficiently low to enable real-time security decision-making. Research on efficient RL architectures, transfer learning approaches that leverage knowledge from related domains, and meta-learning techniques that enable rapid adaptation may help address these computational challenges and enable broader deployment.

Future research directions include development of safe exploration methods that constrain learning to acceptable risk levels during training, preventing exploratory actions that could create vulnerabilities or disrupt operations. Integration of domain knowledge and threat intelligence into RL frameworks through structured representations may accelerate learning and improve policy quality by incorporating human expertise. Federated learning approaches enabling collaborative training across multiple organizations while preserving data privacy could address data scarcity challenges and enable development of more robust models trained on diverse attack examples. Research on combining RL with other AI techniques including causal reasoning and symbolic AI may enable more interpretable and reliable security systems.

The successful integration of RL into operational cybersecurity requires continued collaboration between machine learning researchers, security practitioners, and organizational stakeholders to ensure that developed systems address real operational needs while respecting practical constraints. While RL offers powerful capabilities for autonomous learning and adaptation, these techniques represent tools that must be deployed thoughtfully within comprehensive security programs that include human expertise, organizational processes, and complementary technical controls. The future of cybersecurity will likely involve increasing adoption of RL-based systems as challenges around interpretability,

robustness, and sample efficiency are progressively addressed through ongoing research and practical experience with operational deployments.

## References

[1] Sarker IH, Kayes ASM, Badsha S, et al. Cybersecurity data science: an overview from machine learning perspective. Journal of Big Data. 2020;7(1):1-29.

[2] Apruzzese G, Colajanni M, Ferretti L, et al. Addressing adversarial attacks against security systems based on machine learning. IEEE International Conference on Cyber Conflict. 2019:264-279.

[3] Yeboah-Ofori A, Islam S, Lee SW, et al. Cyber threat predictive analytics for improving cyber supply chain security. IEEE Access. 2021;9:94318-94337.

[4] Nguyen TT, Reddi VJ. Deep reinforcement learning for cyber security. IEEE Transactions on Neural Networks and Learning Systems. 2021;32(8):3779-3795.

[5] Luong NC, Hoang DT, Gong S, et al. Applications of deep reinforcement learning in communications and sensing: a survey. IEEE Communications Surveys & Tutorials. 2019;21(4):3133-3174.

[6] Dulac-Arnold G, Levine N, Mankowitz DJ, et al. Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. Machine Learning. 2021;110(9):2419-2468.

[7] Haarnoja T, Zhou A, Abbeel P, Levine S. Soft actor-critic: off-policy maximum entropy deep reinforcement learning with a stochastic actor. International Conference on Machine Learning. 2019:1861-1870.

[8] Sethi K, Kumar R, Prajapati N, Bera P. Deep reinforcement learning based intrusion detection system for cloud infrastructure. IEEE International Conference on Communication Systems. 2020:1-6.

[9] Alghamdi, F. (2025). Automated Penetration Testing Through Reinforcement Learning. In Complexities and Challenges for Securing Digital Assets and Infrastructure (pp. 323-352). IGI Global Scientific Publishing.

[10] Usama M, Qadir J, Raza A, et al. Unsupervised machine learning for networking: techniques, applications and research challenges. IEEE Access. 2019;7:65579-65615.

[11] Fathy, M., & Tarek, H. (2024). Developing Effective Incident Response Plans: Maintaining Information Assurance During and After Security Breaches. Journal of Data Science, Predictive Analytics, and Big Data Applications, 9(11), 1-14.

[12] Ahmed M, Seraj R, Islam SMS. The k-means algorithm: a comprehensive survey and performance evaluation. Electronics. 2020;9(8):1295.

[13] Chukwuani, E. N., Odunsi, O. R., & Ikemefuna, C. D. (2025). Machine learning techniques for real-time malware classification and threat detection in distributed systems.

[14] Andrychowicz OM, Baker B, Chociej M, et al. Learning dexterous in-hand manipulation. International Journal of Robotics Research. 2020;39(1):3-20.

[15] Zhang K, Yang Z, Başar T. Multi-agent reinforcement learning: a selective overview of theories and algorithms. Handbook of Reinforcement Learning and Control. 2021:321-384.

[16] Huang L, Zhu Q. Analysis and computation of adaptive defense strategies against advanced persistent threats for cyber-physical systems. ACM Transactions on Modeling and Computer Simulation. 2019;29(4):1-28.

[17] Huang L, Zhu Q. Adaptive strategic cyber defense for advanced persistent threats in critical infrastructure networks. ACM SIGMETRICS Performance Evaluation Review. 2019;46(2):52-56.

[18] Vinyals O, Babuschkin I, Czarnecki WM, et al. Grandmaster level in StarCraft II using multi-agent reinforcement learning. Nature. 2019;575(7782):350-354.

[19] Feng S, Haykin S. Cognitive risk control for anti-jamming V2V communications in autonomous vehicle networks. IEEE Transactions on Vehicular Technology. 2019;68(10):9920-9934.

[20] Brown D, Coleman R, Srinivasan R, Niekum S. Safe imitation learning via fast Bayesian reward inference from preferences. International Conference on Machine Learning. 2020:1165-1177.

[21] Dulac-Arnold G, Mankowitz D, Hester T. Challenges of real-world reinforcement learning. arXiv preprint. 2019;arXiv:1904.12901.

[22] Ilahi I, Usama M, Qadir J, et al. Challenges and countermeasures for adversarial attacks on deep reinforcement learning. IEEE Transactions on Artificial Intelligence. 2022;3(2):90-109.

[23] Ball P, Parker-Holder J, Pacchiano A, et al. Ready policy one: world building through active learning. International Conference on Machine Learning. 2020:714-725.

[24] Arrieta AB, Díaz-Rodríguez N, Del Ser J, et al. Explainable artificial intelligence: concepts, taxonomies, opportunities and challenges toward responsible AI. Information Fusion. 2020;58:82-115.

[25] Rudin C, Chen C, Chen Z, et al. Interpretable machine learning: fundamental principles and 10 grand challenges. Statistics Surveys. 2022;16:1-85.

[26] Guidotti R, Monreale A, Ruggieri S, et al. A survey of methods for explaining black box models. ACM Computing Surveys. 2019;51(5):1-42.

[27] Gautam, M. (2023). Deep Reinforcement learning for resilient power and energy systems: Progress, prospects, and future avenues. Electricity, 4(4), 336-380.

[28] Villar-Martínez, A., Rodriguez-Gil, L., Angulo, I., Orduña, P., García-Zubía, J., & López-De-Ipiña, D. (2019). Improving the scalability and replicability of embedded systems remote laboratories through a cost-effective architecture. IEEE Access, 7, 164164-164185.

[29] Hindy H, Brosset D, Bayne E, et al. A taxonomy of network threats and the effect of current datasets on intrusion detection systems. IEEE Access. 2020;8:104650-104675.

[30] Lopez-Martin M, Carro B, Sanchez-Esguevillas A. Application of deep reinforcement learning to intrusion detection for supervised problems. Expert Systems with Applications. 2020;141:112963.

[31] Han D, Wang Z, Zhong Y, et al. Evaluating and improving adversarial robustness of machine learning-based network intrusion detectors. IEEE Journal on Selected Areas in Communications. 2021;39(8):2632-2647.

[32] Ferrag MA, Maglaras L, Moschoyiannis S, Janicke H. Deep learning for cyber security intrusion detection: approaches, datasets, and comparative study. Journal of Information Security and Applications. 2020;50:102419.

[33] Zhao, D., Liu, J., Wang, J., Niu, W., Tong, E., Chen, T., & Li, G. (2019). Bidirectional RNN-based few-shot training for detecting multi-stage attack. arXiv preprint arXiv:1905.03454.

[34] Xing, S., & Wang, Y. (2025). Cross-Modal Attention Networks for Multi-Modal Anomaly Detection in System Software. *IEEE Open Journal of the Computer Society*.

[35] Ring M, Wunderlich S, Scheuring D, et al. A survey of network-based intrusion detection data sets. Computers & Security. 2019;86:147-169.

[36] Sarker, K. U., Yunus, F., & Deraman, A. (2023). Penetration taxonomy: A systematic review on the penetration process, framework, standards, tools, and scoring methods. Sustainability, 15(13), 10471.

[37] Hu Z, Beuran R, Tan Y. Automated penetration testing using deep reinforcement learning. IEEE European Symposium on Security and Privacy Workshops. 2020:2-10.

[38] Schwartz J, Kurniawati H. Autonomous penetration testing using reinforcement learning. arXiv preprint. 2019;arXiv:1905.05965.

[39] Chaudhary S, O'Brien A, Xu S. Automated post-breach penetration testing through reinforcement learning. IEEE Conference on Communications and Network Security. 2020:1-2.

[40] Chen, T., Liu, J., Xiang, Y., Niu, W., Tong, E., & Han, Z. (2019). Adversarial attack and defense in reinforcement learning-from AI security view. Cybersecurity, 2(1), 11.

[41] Zhang, L., & Zhang, L. (2022). Artificial intelligence for remote sensing data analysis: A review of challenges and opportunities. IEEE Geoscience and Remote Sensing Magazine, 10(2), 270-294.

[42] Ye, J., Giani, A., Elasser, A., Mazumder, S. K., Farnell, C., Mantooth, H. A., ... & Abbaszada, M. A. (2021). A review of cyber–physical security for photovoltaic systems. IEEE Journal of Emerging and Selected Topics in Power Electronics, 10(4), 4879-4901.

[43] Wang, Y., Wang, J., Zhang, W., Zhan, Y., Guo, S., Zheng, Q., & Wang, X. (2022). A survey on deploying mobile deep learning applications: A systemic and technical perspective. Digital Communications and Networks, 8(1), 1-17.

[44] Nguyen TT, Nguyen ND, Nahavandi S. Deep reinforcement learning for multiagent systems: a review of challenges, solutions, and applications. IEEE Transactions on Cybernetics. 2020;50(9):3826-3839.

[45] Morić, Z., Dakić, V., & Regvart, D. (2025). Advancing Cybersecurity with Honeypots and Deception Strategies. In Informatics (Vol. 12, No. 1, p. 14). MDPI AG.

[46] Lowe R, Wu Y, Tamar A, et al. Multi-agent actor-critic for mixed cooperative-competitive environments. Advances in Neural Information Processing Systems. 2017. Updated to: Yu C, Velu A, Vinitsky E, et al. The surprising effectiveness of PPO in cooperative multi-agent games. Advances in Neural Information Processing Systems. 2022;35:24611-24624.

[47] Seraj E, Silva A, Gombolay M. Multi-UAV planning for cooperative wildfire coverage and tracking with quality-of-service guarantees. Autonomous Agents and Multi-Agent Systems. 2022;36(2):1-39.

[48] Eliyan, L. F., & Di Pietro, R. (2021). DoS and DDoS attacks in Software Defined Networks: A survey of existing solutions and research challenges.

[49] Song, Y., Zhang, D., Wang, J., Wang, Y., Wang, Y., & Ding, P. (2025). Application of deep learning in malware detection: a review. Journal of Big Data, 12(1), 99.

[50] Anderson HS, Kharkar A, Filar B, et al. Learning to evade static PE machine learning malware models via reinforcement learning. arXiv preprint. 2018. Updated to: Fang Z, Wang J, Geng J, Kan X. A survey on ransomware: evolution, taxonomy, and defense solutions. IEEE Access. 2019;7:92820-92838.

[51] Fang, Z., Wang, J., Li, B., Wu, S., Zhou, Y., & Huang, H. (2019). Evading anti-malware engines with deep reinforcement learning. IEEE Access, 7, 48867-48879.

[52] Kalejaiye, A. N. (2022). REINFORCEMENT LEARNING-DRIVEN CYBER DEFENSE FRAMEWORKS: AUTONOMOUS DECISION-MAKING FOR DYNAMIC RISK PREDICTION AND ADAPTIVE THREAT RESPONSE STRATEGIES. International Journal of Engineering Technology Research & Management (IJETRM), 6(12), 92-111.

[53] Chowdhary A, Pisharody S, Alshamrani A, Huang D. Dynamic game based security framework in SDN-enabled cloud networking environments. ACM International Workshop on Security in Software Defined Networks. 2020:53-58.

[54] Kiran, B. R., Sobh, I., Talpaert, V., Mannion, P., Al Sallab, A. A., Yogamani, S., & Pérez, P. (2021). Deep reinforcement learning for autonomous driving: A survey. IEEE transactions on intelligent transportation systems, 23(6), 4909-4926.

[55] Zhu R, Niu Z, Wu X, et al. Safe reinforcement learning with stability guarantee for motion planning of autonomous vehicles. IEEE Robotics and Automation Letters. 2020;5(4):5983-5990.

[56] Singh A, Yang L, Hartikainen K, et al. End-to-end robotic reinforcement learning without reward engineering. Robotics: Science and Systems. 2019:1-11.

[57] Muñoz-González L, Sgandurra D, Paudice A, Lupu EC. Efficient attack graph analysis through approximate inference. ACM Transactions on Privacy and Security. 2019;20(3):1-30.

[58] Chowdhary A, Huang D, Mahendran JS, et al. Autonomous security analysis and penetration testing. IEEE International Conference on Mobile Cloud Computing. 2020:1-10.

[59] Hutter, M., Quarel, D., & Catt, E. (2024). An introduction to universal artificial intelligence. Chapman and Hall/CRC.

[60] Zhu Z, Lin K, Zhou J. Transfer learning in deep reinforcement learning: a survey. IEEE Transactions on Pattern Analysis and Machine Intelligence. 2023;45(4):5149-5165.

[61] Da Silva FL, Costa AHR. A survey on transfer learning for multiagent reinforcement learning systems. Journal of Artificial Intelligence Research. 2019;64:645-703.

[62] Légère, A., Li, L., Rivest, F., & Al Mallah, R. (2024, July). Training Environments for Reinforcement Learning Cybersecurity Agents. In 2024 International Conference on Computing, Internet of Things and Microwave Systems (ICCIMS) (pp. 1-5). IEEE.

[63] Rakelly K, Zhou A, Finn C, et al. Efficient off-policy meta-reinforcement learning via probabilistic context variables. International Conference on Machine Learning. 2019:5331-5340.

[64] Li, Q., Gao, M., Zhang, G., Zhai, W., Chen, J., & Jeon, G. (2024). Towards multimodal disinformation detection by vision-language knowledge interaction. Information Fusion, 102, 102037.

[65] Silver, D., Singh, S., Precup, D., & Sutton, R. S. (2021). Reward is enough. Artificial intelligence, 299, 103535.

[66] Ji S, Pan S, Cambria E, et al. A survey on knowledge graphs: representation, acquisition, and applications. IEEE Transactions on Neural Networks and Learning Systems. 2022;33(2):494-514.

[67] Muneer, A., Waqas, M., Saad, M. B., Showkatian, E., Bandyopadhyay, R., Xu, H., ... & Wu, J. (2025). From Classical Machine Learning to Emerging Foundation Models: Review on Multimodal Data Integration for Cancer Research. arXiv preprint arXiv:2507.09028.

[68] Heuillet A, Couthouis F, Díaz-Rodríguez N. Explainability in deep reinforcement learning. Knowledge-Based Systems. 2021;214:106685.

[69] Rigotti, M., Miksovic, C., Giurgiu, I., Gschwind, T., & Scotton, P. (2021, May). Attention-based interpretability with concept transformers. In International conference on learning representations.

[70] Samek W, Montavon G, Lapuschkin S, et al. Explaining deep neural networks and beyond: a review of methods and applications. Proceedings of the IEEE. 2021;109(3):247-278.