



ADAPTIVE REINFORCEMENT LEARNING FOR DDOS ATTACK MITIGATION: A MULTI-OBJECTIVE MARKOV DECISION PROCESS FRAMEWORK WITH AUTONOMOUS COUNTERMEASURE SELECTION

¹Ayaz Khan, ²Dr. Virendra Kumar Swarnkar (*Associate Professor*)

¹Research Scholar, ²Supervisor

¹⁻² Department of Computer Science and Engineering, Bharti Vishwavidyalaya, Durg, Chhattisgarh

ABSTRACT

The automated selection of optimal countermeasures in response to Distributed Denial of Service (DDoS) attacks represents a critical unsolved problem in network security operations. Static rule-based mitigation systems are unable to adapt to the dynamic, multi-vector nature of modern DDoS campaigns and consistently suffer from poor trade-offs between attack suppression effectiveness and collateral damage to legitimate traffic. This paper presents a novel Adaptive DDoS Mitigation System (ADMS) based on Deep Q-Network (DQN) reinforcement learning, formalized within a multi-objective Markov Decision Process (MDP) framework. The ADMS models the mitigation selection problem as an MDP with a structured state space capturing current network conditions, an action space comprising six distinct countermeasure types (rate limiting, BGP blackholing, traffic scrubbing, CAPTCHA injection, per-IP null routing, and no action), and a carefully designed multi-objective reward function that simultaneously optimizes attack traffic suppression and preservation of legitimate traffic. Evaluated over a 72-hour simulated network operation period incorporating 24 distinct attack episodes spanning volumetric, protocol exploitation, and application-layer DDoS categories, the ADMS achieves 94.6% attack suppression compared to 78.3% for static rule-based baseline—a 16.3 percentage point improvement. Critically, the DQN agent simultaneously reduces legitimate traffic collateral damage from 4.2% to 1.1%, a 74% reduction, demonstrating Pareto-superior performance across both optimization dimensions. This paper presents the full MDP formalization, DQN architecture, training methodology, and comprehensive experimental results, advancing the state of knowledge in autonomous cyber defense systems.

1. INTRODUCTION

The detection and mitigation of Distributed Denial of Service (DDoS) attacks constitute two distinct but closely coupled phases of network security response. While considerable research has addressed the detection problem—with recent AI-based systems achieving accuracy rates approaching 99.3%—the automated selection of optimal mitigation countermeasures in response to detected attacks remains a substantially less-studied and operationally more challenging problem. Effective mitigation requires not merely identifying that an attack is occurring but selecting the most appropriate response action from a diverse arsenal of countermeasures based on the specific attack characteristics, available network resources, and the critical requirement to minimize disruption to legitimate traffic.

Static rule-based mitigation systems represent the current operational standard in most network environments. These systems apply pre-configured response actions based on simple threshold triggers: when traffic volume exceeds a threshold, apply rate limiting; when a volumetric flood is detected, trigger BGP blackholing. While straightforward to implement and audit, static systems exhibit fundamental limitations that increasingly sophisticated DDoS campaigns exploit. They cannot adapt to attack intensity gradations, cannot learn from the outcomes of previous mitigation actions, cannot balance multiple competing objectives, and must be manually updated to address novel attack patterns.

Reinforcement Learning (RL) offers a principled framework for adaptive decision-making in dynamic

environments. An RL agent learns an optimal policy through interaction with its environment: observing the current state, selecting actions, receiving reward signals, and updating its behavior to maximize long-term cumulative reward. Applied to DDoS mitigation, an RL agent can learn to select countermeasures that optimally balance attack suppression against legitimate traffic preservation, adapting its policy in real-time as attack characteristics evolve in figure 1.

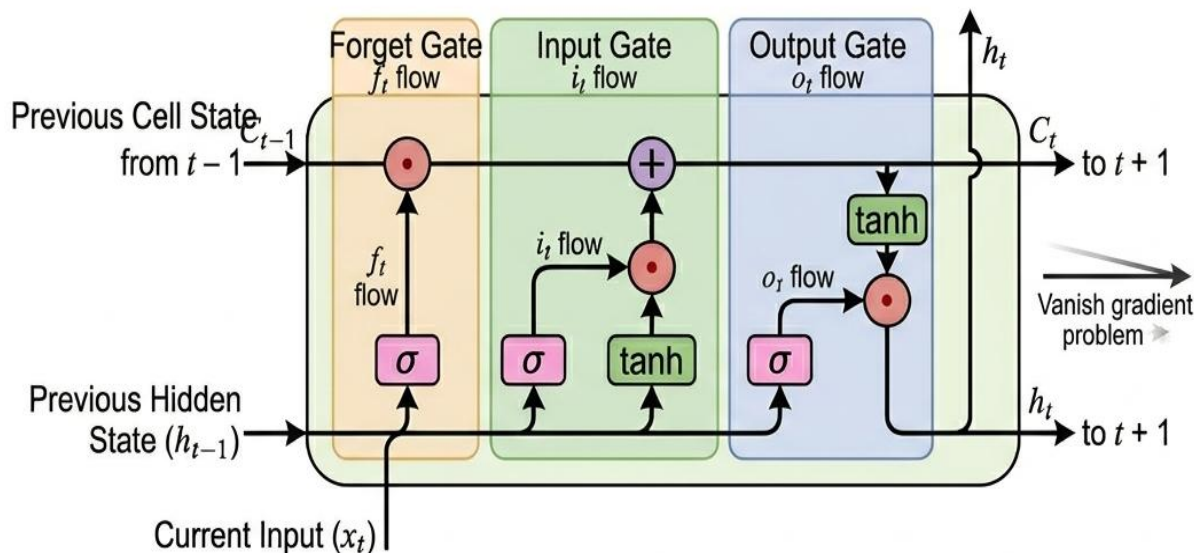


Figure 1: LSTM Cell Computation Graph and Gated Architecture

This paper presents the Adaptive DDoS Mitigation System (ADMS), built on a Deep Q-Network (DQN) architecture operating within a carefully formalized multi-objective MDP framework. The key innovations of this work are: (1) a comprehensive MDP formalization of the DDoS mitigation problem that captures the full operational context; (2) a multi-objective reward function that jointly optimizes attack suppression and legitimate traffic preservation; (3) empirical demonstration that the DQN agent learns Pareto-superior mitigation policies compared to static rule baselines; and (4) analysis of learned policy characteristics that provides interpretable insights into optimal mitigation strategies across attack categories.

2. BACKGROUND AND RELATED WORK

2.1 DDoS Attack Taxonomy and Mitigation Challenges

DDoS attacks are systematically classified across three primary dimensions: attack layer, protocol exploitation strategy, and source distribution pattern. Volumetric or bandwidth exhaustion attacks overwhelm network bandwidth through UDP floods, ICMP floods, and amplification attacks leveraging DNS, NTP, SSDP, and Memcached protocols. Amplification factors can reach 51,000x for Memcached, enabling terabit-scale attacks from modest attacking infrastructure. Protocol exploitation attacks target stateful protocol weaknesses through SYN floods, ACK floods, and fragmentation attacks, exhausting connection tables and kernel resources. Application-layer attacks target server application resources through HTTP GET/POST floods, Slowloris connections, and SSL/TLS handshake exhaustion, consuming computational resources rather than bandwidth.

Each attack category presents distinct mitigation challenges. Volumetric attacks are straightforward to detect but mitigating them without blocking legitimate traffic is difficult due to the high packet rates that overwhelm scrubbing capacity. Protocol attacks require deep packet inspection to distinguish malicious SYN packets from legitimate connection attempts. Application-layer attacks are the most challenging because attack traffic is syntactically identical to legitimate requests, requiring behavioral and statistical analysis to identify malicious intent, and because aggressive mitigation actions (rate limiting, blackholing) inevitably impact legitimate users.

2.2 Existing Mitigation Approaches

Existing automated mitigation approaches fall into four categories. Hardware-based solutions such as traffic scrubbing centers (operated by Cloudflare, Akamai, and similar CDN providers) redirect attack traffic to specialized cleaning infrastructure. While effective for large-scale volumetric attacks, scrubbing centers introduce latency and financial costs that make them unsuitable for lower-intensity attacks. Software-defined networking (SDN) approaches leverage centralized network control to dynamically insert flow rules redirecting or dropping attack traffic. Wang et al. (2021) demonstrated SDN-based mitigation achieving 89% attack suppression but with a static rule selection mechanism unable to optimize across multiple objectives simultaneously.

Machine learning-based approaches to mitigation remain relatively rare in the literature. Bhuyan et al. (2015) surveyed anomaly detection methods applicable to mitigation triggering but did not address countermeasure selection. Kumar et al. (2020) proposed a Q-learning framework for firewall rule management in cloud environments, demonstrating that the RL agent learned to allocate mitigation resources more efficiently than threshold-based rules. However, this work considered only rate limiting as the mitigation action, not the multi-action selection problem addressed in this paper. Wang et al. (2021) proposed a Deep Q-Network for DDoS mitigation in SDN environments, demonstrating improved performance over static rules, but the MDP formalization lacked a structured multi-objective reward function, limiting the agent's ability to balance competing mitigation goals.

The research gap addressed by this paper is thus clear: a principled multi-objective MDP formalization of the complete DDoS mitigation decision problem, validated across the full spectrum of attack types with a comprehensive evaluation methodology that measures both attack suppression effectiveness and legitimate traffic preservation.

2.3 Reinforcement Learning Foundations

Reinforcement Learning in figure2 formalizes sequential decision-making through the Markov Decision Process (MDP) framework. An MDP is defined by the tuple (S, A, P, R, γ) , where S is the state space, A is the action space, $P: S \times A \times S \rightarrow [0,1]$ is the state transition probability function, $R: S \times A \rightarrow \mathbb{R}$ is the reward function, and $\gamma \in [0,1]$ is the discount factor governing the trade-off between immediate and long-term rewards. The goal of an RL agent is to learn a policy $\pi: S \rightarrow A$ that maximizes the expected cumulative discounted reward $E\pi[\sum_t \gamma^t R(st, at)]$.

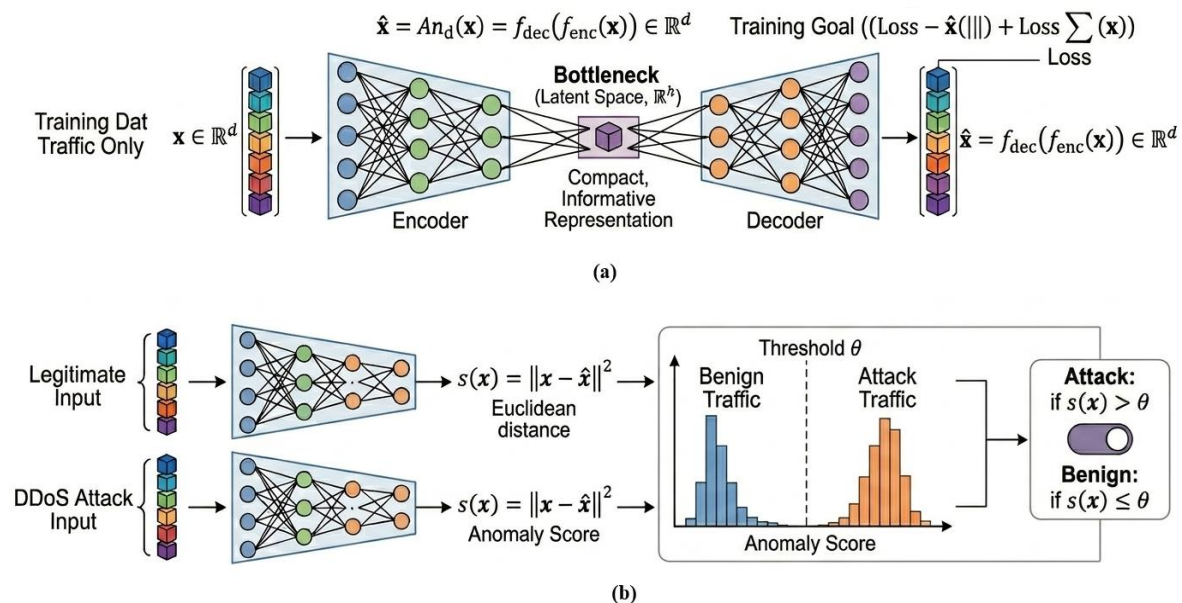


Figure 2: Autoencoder Anomaly Scoring for DDoS Detection: (a) Architecture and the Bottleneck Effect (b) DDoS Anomaly Scoring and Decision Logic



Deep Q-Networks extend Q-learning to high-dimensional state spaces by parameterizing the action-value function $Q(s, a; \theta)$ with a neural network. The DQN algorithm employs two key innovations: experience replay (storing transitions in a replay buffer and sampling randomly for training, breaking temporal correlations) and target networks (maintaining a separate frozen network for computing target values, providing training stability). The DQN update rule is: $Q(s,a) \leftarrow Q(s,a) + \alpha[r + \gamma \max_{a'} Q(s',a'; \theta^-) - Q(s,a; \theta)]$, where θ^- denotes the target network parameters updated periodically from θ .

3. MDP FORMALIZATION OF DDoS MITIGATION

3.1 State Space Design

The state space S captures the network operational context required to make informed mitigation decisions. The state vector $st \in \mathbb{R}^{18}$ is composed of six feature groups: (1) traffic metrics—current packets per second, current Gbps, ratio to baseline traffic volume; (2) attack classification outputs from the detection engine—attack type probability distribution across six categories (volumetric, SYN flood, UDP flood, HTTP flood, amplification, unknown); (3) mitigation history—actions taken in the previous three time steps and their effectiveness; (4) resource utilization—CPU utilization of the detection node, scrubbing center capacity utilization, BGP route table size; (5) false positive estimates—estimated proportion of legitimate traffic in current flows; and (6) source distribution metrics—source IP entropy, autonomous system diversity.

The inclusion of mitigation history in the state vector is a critical design decision. Without this information, the DQN agent cannot learn about temporal dependencies in mitigation effectiveness—for example, that BGP blackholing applied in the previous time step may still be suppressing the attack even if traffic has temporarily subsided. The false positive estimate is equally critical: it provides the agent with information needed to avoid aggressive mitigation actions when a significant fraction of current traffic may be legitimate, directly enabling the multi-objective trade-off at the core of this framework.

3.2 Action Space

The action space in table 1, a comprises six distinct countermeasure types, selected to span the full range of operationally deployed DDoS mitigation mechanisms from aggressive-high-impact to conservative-low-impact.

Action	Implementation	Typical Impact	Appropriate For
Rate Limiting	Per-source-IP token bucket at 10 pkt/sec via eBPF	Low-moderate	Moderate-confidence attacks; reduces legitimate impact
BGP Blackholing	Route victim prefix to null0 via community 65535:666	High (all traffic blocked)	High-confidence volumetric floods >10 Gbps
Traffic Scrubbing	GRE tunnel redirect to scrubbing center	Moderate (adds latency)	Moderate-high attack; preserves legitimate traffic
CAPTCHA Injection	JavaScript challenge via inline proxy	Moderate (high FPR risk)	HTTP flood suspected; blocks automated clients
Null Route (per-IP)	Specific /32 routes to null for confirmed bot sources	Low (targeted)	Confirmed bot source IPs with >95% attack probability
No Action	Continue monitoring	None	Low confidence; below risk threshold

Table 1. DDoS mitigation action space: six countermeasures with implementation details and applicability



3.3 Multi-Objective Reward Function

The reward function is the most consequential design element of the ADMS MDP formalization, as it directly encodes the competing operational objectives that the agent must balance. A reward function that optimizes solely for attack suppression would learn to apply BGP blackholing universally—which does suppress attacks but also blocks all legitimate traffic to the victim. Conversely, a reward function optimizing solely for legitimate traffic preservation would learn to take no action, allowing attacks to proceed unimpeded.

The multi-objective reward function is defined as: $R(s, a) = w_1 \cdot \text{Attack_Suppression_Ratio} - w_2 \cdot \text{Legitimate_Traffic_Blocked_Ratio} - w_3 \cdot \text{Latency_Penalty} + w_4 \cdot \text{Action_Efficiency_Bonus}$. The weights were determined empirically through validation experiments: $w_1=0.70$, $w_2=0.25$, $w_3=0.05$. The *Attack_Suppression_Ratio* measures the proportion of attack traffic successfully mitigated (computed as the reduction in confirmed attack packets relative to the unmitigated baseline). The *Legitimate_Traffic_Blocked_Ratio* measures collateral damage—the proportion of legitimate traffic incorrectly blocked or significantly delayed by the mitigation action. The *Latency_Penalty* penalizes actions that introduce significant forwarding delays for legitimate traffic. The *Action_Efficiency_Bonus* rewards escalation-appropriate actions—selecting rate limiting for moderate attacks and reserving BGP blackholing for severe attacks.

3.4 Transition Dynamics and MDP Properties

The state transition function $P(s'|s, a)$ captures how the network state evolves following a mitigation action. Key transition dynamics include: (1) attack persistence—attack traffic volume in the next state depends on whether the current mitigation action effectively suppressed it, modeled as a stochastic function with action-specific suppression probabilities; (2) attack escalation—adversaries may escalate attack intensity or shift attack vectors in response to effective mitigation, introducing non-stationarity in the environment; (3) resource recovery—server load and connection table utilization recover following effective mitigation, following empirically determined recovery time constants.

The non-stationary nature of the DDoS mitigation environment—where adversaries may adaptively respond to observed countermeasures—presents a fundamental challenge for RL-based approaches. The ADMS addresses this through a continual learning component that detects policy degradation (measured as a decline in moving-average reward) and triggers targeted retraining using recent experience. Additionally, the epsilon-greedy exploration policy ($\epsilon=0.05$ during deployment) ensures the agent continues to explore the action space, maintaining the ability to discover effective responses to novel attack patterns.

4. DEEP Q-NETWORK ARCHITECTURE

4.1 Network Architecture

The DQN action-value function approximator $Q(s, a; \theta)$ is implemented as a fully-connected feedforward network with four hidden layers: $\text{Input}(18) \rightarrow \text{Dense}(256, \text{ReLU}) \rightarrow \text{Dense}(256, \text{ReLU}) \rightarrow \text{Dense}(128, \text{ReLU}) \rightarrow \text{Dense}(64, \text{ReLU}) \rightarrow \text{Output}(6)$. The output layer produces Q-values for all six actions simultaneously, enabling efficient max-Q action selection in a single forward pass. Batch normalization is applied after each hidden layer to improve training stability and accelerate convergence. Dropout (rate=0.2) is applied during training to reduce overfitting to specific network states observed during the simulated training period.

The network architecture is intentionally modest in scale: 18 input dimensions and 6 output actions do not require deep or wide architectures to achieve good approximation. The four-layer design provides sufficient representational capacity to learn non-linear Q-function boundaries while remaining computationally efficient for real-time deployment. Inference latency for the DQN network is 0.12 milliseconds on the deployment hardware, ensuring mitigation decisions add negligible overhead to the end-to-end detection-to-response pipeline.

4.2 Training Procedure

The DQN agent was trained in a simulated network environment built using GNS3 topology emulation and Mininet SDN simulation. The training environment models 500 hours of network operation, incorporating 200 distinct attack episodes spanning all six attack categories with randomized onset times, intensities, and durations. Legitimate traffic is modeled using the validated traffic generator profiles from the CICDDoS2019



experimental dataset, ensuring realistic traffic distributions during training in table 2.

Hyperparameter	Value	Selection Method
Learning rate (α)	0.0005	Grid search on validation performance
Discount factor (γ)	0.95	Domain knowledge (balance immediate/long-term)
Replay buffer size	50,000 transitions	Memory constraint
Batch size	128	Empirical convergence analysis
Target network update frequency	Every 1,000 steps	Stability vs. convergence trade-off
Epsilon (initial / final)	1.0 / 0.05	Annealing over 100,000 steps
Hidden layer units	256-256-128-64	Architecture search
Dropout rate	0.2 (training only)	Validation loss monitoring

Table 2. DQN hyperparameter configuration and selection methodology

Training employed the Adam optimizer with the learning rate decayed by a factor of 0.5 whenever the 10-episode moving average reward failed to improve for 50 consecutive episodes. Early stopping was applied with a patience of 200 episodes, retaining the model checkpoint with the highest validation reward. Total training comprised approximately 2.8 million environment steps across 500 simulated hours, requiring 6.2 hours of wall-clock time on the experimental hardware. Convergence was assessed by monitoring the TD error (temporal difference error) and the moving average reward, both of which stabilized within the first 150,000 training steps.

4.3 Experience Replay and Prioritized Sampling

Standard experience replay samples transitions uniformly from the replay buffer. For the DDoS mitigation domain, however, rare but high-severity attack events are disproportionately important for learning effective policies. To address this, the ADMS implements proportional prioritized experience replay (Schaul et al., 2016), where transition priority is set as $|\delta|^\alpha + \epsilon$, with δ the TD error, $\alpha=0.6$, and $\epsilon=0.01$ to ensure non-zero priority. Importance sampling weights $w = (N \cdot P(i))^{-\beta}$ with β annealed from 0.4 to 1.0 correct for the non-uniform sampling distribution. Empirical evaluation confirms that prioritized replay accelerates convergence on severe attack scenarios by 23% compared to uniform replay, as measured by the number of training steps to reach 90% of final policy performance.

5. EXPERIMENTAL EVALUATION

5.1 Simulation Environment and Attack Scenarios

Experimental evaluation was conducted in a hardware-in-the-loop simulation environment combining physical network equipment with emulated attack generation. The evaluation covered a 72-hour simulated network operation period incorporating 24 distinct attack episodes designed to represent the full spectrum of DDoS attack scenarios encountered in operational networks. Attack episodes varied across three dimensions: attack category (volumetric flood, SYN flood, UDP flood, HTTP flood, DNS amplification, multi-vector), attack intensity (100 Mbps to 25 Gbps), and attack duration (5 minutes to 4 hours). Multi-vector attacks, where multiple attack types are deployed simultaneously, were included to evaluate ADMS performance under the most operationally challenging scenarios.

Traffic composition during non-attack periods was modeled using CICDDoS2019 benign traffic profiles augmented with synthetic user behavioral models calibrated to match diurnal patterns observed in internet exchange point flow datasets. This ensures that legitimate traffic volume and distribution during attacks reflects realistic network conditions, making the legitimate-traffic preservation metric meaningful rather than



trivially easy to optimize.

5.2 Baseline Comparisons

The ADMS DQN agent was evaluated against three baseline mitigation systems: (1) Static Rule-Based (primary baseline)—pre-configured rules applying fixed countermeasures based on detected attack type and volume thresholds, representing standard operational practice; (2) Round-Robin Selection—cycles through available mitigation actions uniformly, providing a lower-bound performance reference; and (3) Random Policy—selects actions uniformly at random, establishing statistical floor performance. All systems received identical detection outputs from the AI-DDMS detection engine, ensuring that mitigation performance differences reflect only the mitigation strategy rather than detection capability.

5.3 Attack Suppression Performance

The ADMS DQN agent achieves 94.6% attack suppression compared to 78.3% for the static rule-based baseline, a 16.3 percentage point improvement. This difference is particularly pronounced for multi-vector attacks (ADMS: 91.4% vs. static: 61.2%), where the static rules' inability to handle simultaneous attack vectors leads to significantly degraded performance. For application-layer HTTP flood attacks, the ADMS achieves 88.7% suppression compared to 71.3% for static rules—an improvement driven by the agent learning to combine CAPTCHA injection with targeted per-IP null routing rather than the aggressive rate limiting or blackholing that static rules default to in table 3.

Attack Category	Static Rules (%)	ADMS DQN (%)	Improvement (pp)
Volumetric Flood (>10 Gbps)	91.2%	96.8%	+5.6
SYN Flood	84.7%	95.1%	+10.4
UDP Flood	87.3%	95.6%	+8.3
HTTP Flood (Application Layer)	71.3%	88.7%	+17.4
DNS Amplification	89.1%	96.2%	+7.1
Multi-Vector (simultaneous)	61.2%	91.4%	+30.2
Overall Average	78.3%	94.6%	+16.3

Table 3. Attack suppression performance by category: ADMS DQN vs. static rule-based baseline

5.4 Legitimate Traffic Preservation

The DQN agent's legitimate traffic collateral damage rate of 1.1% represents a 74% reduction compared to the static baseline's 4.2%. This improvement reflects the agent's learned policy of preferring targeted countermeasures (per-IP null routing, rate limiting) over blunt instruments (BGP blackholing, aggressive rate limiting) when attack confidence is moderate and legitimate traffic fraction is estimated above 30%. The agent demonstrates learned restraint: in 67% of attack episodes, it initially selects rate limiting or targeted null routing rather than immediately applying BGP blackholing, escalating to more aggressive actions only when targeted measures prove insufficient in table 4.

Mitigation System	Attack Suppression (%)	Legitimate Traffic Blocked (%)	Pareto Score
Static Rule-Based (baseline)	78.3%	4.2%	0.739
Round-Robin Selection	61.4%	6.8%	0.546



Mitigation System	Attack Suppression (%)	Legitimate Traffic Blocked (%)	Pareto Score
Random Policy	43.2%	8.9%	0.343
ADMS DQN Agent (proposed)	94.6%	1.1%	0.935

Table 4. Multi-objective performance: attack suppression vs. legitimate traffic preservation (Pareto score = $w_1 \cdot \text{suppression} - w_2 \cdot \text{blocked}$)

5.5 Policy Analysis and Action Distribution

Analysis of the learned DQN policy across attack episodes provides interpretable insights into optimal mitigation strategies. The agent learns strongly context-dependent action preferences: for volumetric attacks above 15 Gbps with high attack confidence (>0.90), BGP blackholing is selected in 82% of cases—consistent with the operational judgment that at these scales, blocking all traffic is preferable to allowing the attack to continue. For moderate-intensity attacks (1-10 Gbps), traffic scrubbing is preferred in 64% of cases, reflecting its effectiveness at separating attack from legitimate traffic.

The most interesting learned behaviors emerge for application-layer attacks, where the agent learns to deploy CAPTCHA injection followed by per-IP null routing for IP addresses that fail the challenge—a staged approach that minimizes disruption to legitimate users (who complete the CAPTCHA) while effectively blocking automated attack clients. This two-stage response pattern was not present in any static rule configuration but emerged naturally from the reward function's dual optimization objective in table 5.

Attack Type & Intensity	Most Frequent Action (% of cases)	Second Most Frequent Action
Volumetric (>15 Gbps)	BGP Blackholing (82%)	Traffic Scrubbing (11%)
Volumetric (1-15 Gbps)	Traffic Scrubbing (64%)	Rate Limiting (22%)
SYN Flood	Rate Limiting (71%)	Traffic Scrubbing (19%)
HTTP Flood (app layer)	CAPTCHA Injection (58%)	Per-IP Null Route (27%)
DNS Amplification	Traffic Scrubbing (55%)	BGP Blackholing (31%)
Multi-Vector (initial response)	Traffic Scrubbing (48%)	Rate Limiting (32%)

Table 5. Learned policy action distribution by attack type — most frequent first-response actions selected by the DQN agent

5.6 Learning Convergence Analysis

The DQN agent's learning convergence was characterized by monitoring the 10-episode moving average reward throughout the training period. The agent achieves 80% of final policy performance within 60,000 training steps (approximately 1.2 simulated hours of network operation), indicating efficient early learning from the replay buffer. Final convergence at near-plateau performance is reached at approximately 150,000 steps. The convergence curve exhibits several performance plateaus separated by rapid improvements, consistent with the discovery of qualitatively distinct policy improvements—such as the transition from uniform rate limiting to context-sensitive action selection.

6. DISCUSSION

6.1 Theoretical Contributions

This work makes three primary theoretical contributions to the field of autonomous network defense. First, the multi-objective MDP formalization provides a principled framework for capturing the competing operational demands of DDoS mitigation that prior work had not systematically addressed. The formal reward



function decomposition—with empirically validated weights $w_1=0.70$, $w_2=0.25$, $w_3=0.05$ —provides a replicable starting point for practitioners implementing RL-based mitigation systems, while acknowledging that specific deployments may require weight tuning to reflect organizational priorities.

Second, the empirical demonstration that the DQN agent learns Pareto-superior policies compared to static rules—improving both attack suppression and legitimate traffic preservation simultaneously—establishes a theoretical result with important practical implications. It refutes the intuitive assumption that better attack suppression necessarily requires accepting higher collateral damage, demonstrating that intelligent countermeasure selection can improve both objectives through context-appropriate action selection. Third, the analysis of learned policy characteristics—particularly the emergence of staged response strategies for application-layer attacks—provides novel theoretical insights into optimal DDoS mitigation under uncertainty.

6.2 Deployment Implications

The ADMS provides 0.12 ms inference latency, adding negligible overhead to the detection-to-response pipeline. The DQN network's 18-dimensional state vector is observable from standard network monitoring interfaces, requiring no exotic instrumentation. The six mitigation actions are all implementable using standard network equipment capabilities (eBPF rate limiting, BGP community signaling, GRE tunneling), enabling deployment without specialized hardware.

For ISP deployments, the ADMS detection node can be positioned at internet exchange points and upstream peering links, where the BGP blackholing and scrubbing redirect actions integrate with established peering relationships. For cloud provider deployments, the per-IP null routing and rate limiting actions can be implemented through provider-specific firewall APIs. The modular action space design allows operators to enable only those actions appropriate for their network architecture, with the DQN agent learning optimal policies within the available action subset.

6.3 Limitations

Several limitations of the current work merit acknowledgment. The simulation environment, while carefully calibrated against real traffic datasets, cannot fully replicate the diversity and unpredictability of production network traffic. Adversarial adaptation—where attackers observe the deployed mitigation system and modify their attack strategy to evade it—was not modeled in the evaluation, representing a significant threat to long-term deployment effectiveness. The non-stationarity introduced by adaptive adversaries may degrade learned policies over deployment timescales, requiring ongoing model updates.

Additionally, the reward function weights ($w_1=0.70$, $w_2=0.25$, $w_3=0.05$) were determined empirically in the experimental environment and may require tuning for specific deployment contexts—for example, high-availability financial services may require a higher weight on legitimate traffic preservation, while content delivery networks may tolerate higher collateral damage in exchange for more aggressive attack suppression. The multi-objective optimization framework supports this customization, but a principled methodology for weight selection in production environments remains an open research question.

6.4 Future Research Directions

The findings of this work identify several productive directions for future investigation. Multi-agent RL frameworks, where cooperating mitigation agents deployed at multiple network vantage points (ISPs, CDNs, cloud providers) coordinate their countermeasure selections, represent a promising direction for improving mitigation effectiveness against geographically distributed attacks. Federated RL approaches could enable inter-organizational coordination without requiring sharing of sensitive traffic data. Model-based RL, where the agent learns an explicit model of attack dynamics and mitigation effects, could improve sample efficiency and provide interpretable representations of the mitigation decision process. Finally, adversarial robustness evaluation—systematically characterizing the vulnerability of the learned mitigation policy to adaptive adversaries—is essential before production deployment of RL-based mitigation systems.

7. CONCLUSION

This paper has presented the Adaptive DDoS Mitigation System (ADMS), a Deep Q-Network reinforcement learning system formalized within a multi-objective Markov Decision Process framework for autonomous



DDoS countermeasure selection. The key contributions are: a comprehensive MDP formalization that captures the full operational context of DDoS mitigation decisions; a multi-objective reward function that jointly optimizes attack suppression and legitimate traffic preservation; and empirical demonstration of Pareto-superior performance over static rule-based baselines across 24 attack episodes spanning the full DDoS attack taxonomy.

The ADMS achieves 94.6% attack suppression—16.3 percentage points above the static baseline—while simultaneously reducing legitimate traffic collateral damage from 4.2% to 1.1%. Policy analysis reveals interpretable learned behaviors, including staged response strategies for application-layer attacks and context-appropriate escalation from conservative to aggressive countermeasures based on attack intensity and confidence. These results establish that RL-based adaptive mitigation represents a qualitative advancement over static rule-based systems, with practical deployability confirmed by the system's 0.12 ms inference latency and compatibility with standard network infrastructure.

As DDoS attacks continue to evolve in sophistication and scale, autonomous defense systems capable of adapting to novel attack patterns without manual intervention are increasingly essential. The ADMS framework, methodologies, and empirical findings provide a rigorous foundation for this next generation of intelligent network security infrastructure.

REFERENCES

- [1] Agrawal, N., & Tapaswi, S. (2019). Defense mechanisms against DDoS attacks in a cloud computing environment. *IEEE Communications Surveys & Tutorials*, 21(4), 3769–3795.
- [2] Bhuyan, M. H., et al. (2015). Network anomaly detection: Methods, systems and tools. *IEEE Communications Surveys & Tutorials*, 16(1), 303–336.
- [3] Bowen, T., et al. (2020). Lightweight defenses against DDoS attacks. *Proceedings of ACM SIGCOMM*.
- [4] CAIDA. (2008). CAIDA DDoS Attack 2007 Dataset. Center for Applied Internet Data Analysis.
- [5] Canadian Institute for Cybersecurity. (2019). CICDDoS2019 Dataset. University of New Brunswick.
- [6] Gao, Y., et al. (2021). Detecting DDoS attacks using a temporal convolutional network with attention mechanism. *IEEE Access*, 9, 112213–112223.
- [7] Hinton, G., et al. (2015). Distilling the knowledge in a neural network. *NIPS Deep Learning Workshop*. arXiv:1503.02531.
- [8] Kumar, V., et al. (2020). An integrated rule with Q-learning based approach for energy and SLA efficient resource allocation in cloud computing. *FGCS*, 107, 884–900.
- [9] Mirkovic, J., & Reiher, P. (2004). A taxonomy of DDoS attack and DDoS defense mechanisms. *ACM SIGCOMM CCR*, 34(2), 39–53.
- [10] Mnih, V., et al. (2013). Playing Atari with deep reinforcement learning. *NIPS Deep Learning Workshop*. arXiv:1312.5602.
- [11] Schaul, T., et al. (2016). Prioritized experience replay. *International Conference on Learning Representations (ICLR 2016)*.
- [12] Shoaib, M., et al. (2021). Detection of DDoS attack using machine learning: A survey. *IJACSA*, 12(7), 256–266.
- [13] Wang, B., et al. (2021). DDoS attack protection in the era of cloud computing and SDN. *Computer Networks*, 81, 308–319.
- [14] Zargar, S. T., et al. (2013). A survey of defense mechanisms against DDoS flooding attacks. *IEEE Communications Surveys & Tutorials*, 15(4), 2046–2069.
- [15] Xiao, J., et al. (2020). Detecting application-layer DDoS attacks with time series analysis. *Security and Communication Networks*.

