



YAHIA FARES UNIVERSITY OF MEDEA

FACULTY OF SCIENCES

DEPARTMENT OF MATHEMATICS AND COMPUTER SCIENCE

Detection of DoS Attacks Using Reinforcement Learning

Prepared and presented by:

Haitham Attab

Amine Takdenti

Supervisor:

Prof. Brahim Sahmadi

A report submitted in partial fulfillment of the requirements for the degree of Master of
Computer Science

June 20, 2025

To my family, for their endless love and support.

Acknowledgements

I would like to express my sincere gratitude to my supervisor, Professor Brahim Sahmadi, for their invaluable guidance and unwavering support throughout this research.

I also wish to thank [mention other people, e.g., colleagues, friends, family] for their encouragement and insightful discussions.

This work was made possible by [mention any funding or support].

Abstract

With the growing frequency and sophistication of Denial-of-Service (DoS) attacks, robust and adaptive detection mechanisms have become critical to maintaining the availability and reliability of computer networks. These attacks can severely disrupt essential services and lead to substantial financial losses, highlighting the need for intelligent and resilient defense strategies.

This thesis introduces a novel intrusion detection system based on Reinforcement Learning (RL) to detect and mitigate DoS attacks. By framing network security as a sequential decision-making problem, the proposed approach enables an RL agent to learn optimal detection policies through continuous interaction with a simulated network environment. The agent incrementally refines its ability to differentiate between benign and malicious traffic by leveraging a reward-based feedback mechanism that reinforces accurate classifications and penalizes misclassifications.

Key traffic features are systematically extracted and provided as state inputs to the RL agent, facilitating effective and dynamic behavioral analysis. This adaptive learning paradigm enhances the system's capability to recognize evolving attack patterns, ultimately contributing to a more proactive and resilient network defense against DoS threats.

ملخص

ضرورية والمتكيفة المتقدمة الكشف آليات أصبحت ،(DoS) الخدمة حجب هجمات وتعقيد وتيرة تزايد مع مالية خسائر في وتسبب الحيوية الخدمات تعطل أن الهجمات لهذه يمكن وموثوقيتها. الشبكات توافر لضمان ومرنة. ذكية دفاع استراتيجيات إلى الحاجة يُبرز مما كبيرة،

لاكتشاف (Reinforcement Learning) المعزز التعلم على يعتمد جديد تسلسل كشف نظام الأطروحة هذه تقدم العميل يتمكن متسلسلة، قرارات اتخاذ كمشكلة الشبكات أمن تصوير خلال من لها. والتصدي DoS هجمات محاكاة. شبكة بيئة مع المستمر التفاعل خلال من المثلث الكشف سياسات تعلم من المعزز التعلم في (Agent) على قائم راجعة تغذية نظام باستخدام والخبيثة السليمة المرور حركة بين التمييز على تدريجياً قدرته بتحسين ويقوم الخاطئة. التصنيفات ويعاقب الصحيحة التصنيفات يعزز المكافآت،

تحليل من يمكنه مما للعميل، حالة كمدخلات وتقديمها منهجي بشكل الأساسية المرور حركة ميزات استخراج يتم أنماط على التعرف على النظام قدرة تعزيز في التكيفية التعليمية المنهجية هذه تسهم وديناميكية. بفعالية السلوك DoS. تهديدات ضد واستباقية مرونة أكثر شبكي دفاع إلى النهاية في يؤدي مما المتغيرة، الهجمات

Résumé

Face à la fréquence et à la complexité croissantes des attaques par déni de service (DoS), la mise en place de mécanismes de détection robustes et adaptatifs est devenue essentielle pour garantir la disponibilité et la fiabilité des réseaux informatiques. Ces attaques peuvent perturber gravement les services critiques et engendrer des pertes financières importantes, soulignant ainsi la nécessité de stratégies de défense intelligentes et résilientes.

Cette thèse propose un système de détection d'intrusion innovant basé sur l'apprentissage par renforcement (Reinforcement Learning) pour détecter et atténuer les attaques DoS. En formulant la sécurité réseau comme un problème de prise de décision séquentielle, notre approche permet à un agent d'apprentissage de découvrir des politiques de détection optimales via une interaction continue avec un environnement réseau simulé. L'agent améliore progressivement sa capacité à distinguer le trafic bénin du trafic malveillant grâce à un système de rétroaction basé sur des récompenses, qui renforce les classifications correctes et pénalise les erreurs.

Les principales caractéristiques du trafic sont extraites de manière systématique et utilisées comme entrées d'état pour l'agent, facilitant une analyse comportementale efficace et dynamique. Ce paradigme d'apprentissage adaptatif améliore la capacité du système à reconnaître les schémas d'attaque en évolution, contribuant ainsi à une défense réseau plus proactive et résiliente contre les menaces DoS.

Contents

Acknowledgements	2
Abstract	3
ملخص	4
Résumé	5
List of Figures	8
List of Tables	9
List of Acronyms	10
1 Intrusion Detection System (IDS)	15
1.1 Security Fundamentals	15
1.1.1 Core Security Concepts	15
1.2 Denial-of-Service (DoS) Attacks	16
1.2.1 Volumetric Attacks	16
1.2.2 Protocol Attacks	17
1.2.3 Application Layer Attacks	17
1.2.4 Distributed Denial-of-Service (DDoS)	18
1.2.5 Network Traffic Analysis	18
1.3 Intrusion Detection Systems (IDS)	18
1.3.1 IDS Architecture	19
1.3.2 IDS Types and Functionality	20
1.4 Classification in Security Contexts	22
1.4.1 Classification Paradigms	22
1.4.2 Machine Learning Techniques	22
2 Reinforcement Learning Applications: A Comprehensive Guide	23
2.1 Introduction to the Reinforcement Learning Approach	23
2.2 Introduction to Deep Q-Learning	23
2.3 Environment Modeling for DoS Detection	24
2.3.1 State Space Design	24
2.3.2 Action Space Definition	25
2.3.3 Reward Function Design	26
2.3.4 Environment Dynamics	26
2.3.5 Challenges in Applying RL to Security Domains	26
2.3.6 Advantages of RL for DoS Detection	27
2.4 Introduction to the Reinforcement Learning Approach	27
2.5 Introduction to Deep Q-Learning	27
2.6 Environment Modeling for DoS Detection	28

2.6.1	State Space Design	29
2.6.2	Action Space Definition	29
2.6.3	Reward Function Design	30
2.6.4	Environment Dynamics	30
2.6.5	Challenges in Applying RL to Security Domains	31
2.6.6	Advantages of RL for DoS Detection	31
3	Designing DQN Architectures for DoS/DDoS Attack Detection	32
3.1	Introduction	32
3.2	Problem Modeling for DoS/DDoS Detection using DQN	32
3.2.1	Environment Definition	32
3.2.2	Data Preprocessing Pipeline	33
3.2.3	Proposed DQN Architectures for Attack Detection	35
3.2.4	Conclusion	39

List of Figures

1.1	Principles of information security	16
1.1	16
1.2	DoS and DDoS Attack Traffic Comparison	17
1.3	Network-based Intrusion Detection System (NIDS) Architecture	19
1.4	Network-based Intrusion Detection System (NIDS) Diagram	20
1.5	Host-based Intrusion Detection System (HIDS) Diagram	21

List of Tables

List of Acronyms

AI	Artificial Intelligence
CIA	Confidentiality, Integrity, and Availability
CIC	Canadian Institute for Cybersecurity
CNN	Convolutional Neural Network
CPU	Central Processing Unit
CSV	Comma-Separated Values
DDoS	Distributed Denial-of-Service
DNS	Domain Name System
DoS	Denial-of-Service
DQN	Deep Q-Network
FTP	File Transfer Protocol
HIDS	Host-based Intrusion Detection System
HTTP	Hypertext Transfer Protocol
HTTPS	Hypertext Transfer Protocol Secure
ICMP	Internet Control Message Protocol
IDS	Intrusion Detection System
IP	Internet Protocol
IoT	Internet of Things
ML	Machine Learning
MLP	Multi-Layer Perceptron
NIDS	Network-based Intrusion Detection System
NTP	Network Time Protocol
PCAP	Packet Capture
RL	Reinforcement Learning
SIEM	Security Information and Event Management
SSDP	Simple Service Discovery Protocol

SVM Support Vector Machine

TCP Transmission Control Protocol

UDP User Datagram Protocol

General Introduction

Context and Motivation

In the contemporary digital era, the reliability and continuous availability of network services form the backbone of global commerce, communication, and critical infrastructure. A primary threat to this stability is the Denial-of-Service (DoS) attack, and its more potent variant, the Distributed Denial-of-Service (DDoS) attack. These malicious activities aim to render online services and network resources unavailable to legitimate users by overwhelming them with an insurmountable volume of traffic or by exploiting protocol vulnerabilities. As attackers devise increasingly sophisticated and large-scale attacks, the need for intelligent, adaptive, and robust defense mechanisms has become more critical than ever.

Traditional Intrusion Detection Systems (IDS) have served as the first line of defense. However, they often rely on static, signature-based methods that, while effective against known threats, are inherently incapable of identifying novel or zero-day attacks. Anomaly-based systems, which attempt to solve this by creating a baseline of "normal" behavior, frequently suffer from high false positive rates, leading to operational friction and alert fatigue. Furthermore, conventional machine learning classifiers, though an improvement, often operate as static models that require periodic retraining and struggle to adapt in real-time to the dynamic and adversarial nature of cybersecurity threats.

This clear gap in defensive capabilities motivates the exploration of more advanced paradigms. Reinforcement Learning (RL), a subfield of machine learning, offers a compelling alternative. Unlike supervised learning, RL focuses on training an autonomous agent to make optimal decisions through interaction with an environment, guided by a system of rewards and penalties. This approach is inherently adaptive, allowing the agent to learn and refine its strategies over time. The potential to create a self-improving, context-aware detection system that can respond dynamically to evolving threats provides the core motivation for this research.

Problem Statement

The central problem addressed by this thesis is the inadequacy of traditional and static machine learning-based Intrusion Detection Systems in providing adaptive, real-time defense against modern, sophisticated DDoS attacks. These existing systems exhibit several key limitations:

- **Lack of Adaptability:** Signature-based systems cannot detect new attack vectors, while static ML models fail to adapt to evolving attack patterns without manual retraining.
- **High False Positive Rates:** Anomaly detection systems often misclassify legitimate traffic spikes (flash crowds) as attacks, disrupting normal operations.
- **Reactive Nature:** Most systems are reactive, identifying attacks only after they have begun to impact the network, rather than learning to predict and preempt them.

This research seeks to answer the following question: **Can Reinforcement Learning, specifically a Deep Q-Network (DQN) approach, be effectively adapted to the supervised classification task of DDoS detection to create a more dynamic, accurate, and autonomous system?**

Objectives

To address the problem statement, this thesis sets out the following primary objectives:

1. To conduct a thorough review of Denial-of-Service attacks, traditional Intrusion Detection Systems, and the foundational principles of Reinforcement Learning to establish a strong theoretical basis.
2. To design and develop a novel detection framework that reframes the DDoS classification task as a sequential decision-making problem, solvable with a Deep Q-Network (DQN) agent.
3. To implement two distinct neural network architectures for the DQN’s Q-function approximator—a Multi-Layer Perceptron (MLP) and a Convolutional Neural Network (CNN)—to explore different feature representation capabilities.
4. To train and rigorously evaluate the proposed models using a modern, comprehensive benchmark dataset (CIC-DDoS2019), measuring performance based on key metrics such as accuracy, precision, recall, and F1-score.
5. To compare the performance of the MLP-based and CNN-based DQN agents to determine the most effective architecture for this task.

Overall Methodology

The research will be conducted through a structured, multi-stage process:

1. **Data Preprocessing:** The CIC-DDoS2019 dataset will be acquired and subjected to a comprehensive preprocessing pipeline. This includes data cleaning, binary label transformation, feature selection using a Random Forest classifier to identify the most salient features, Min-Max normalization to scale the data, and random undersampling to create a balanced dataset for training.
2. **Model Design:** A custom Reinforcement Learning environment will be simulated from the static dataset. A Deep Q-Network (DQN) agent will be designed, defining the state space (from network features), action space (classify as 'Benign' or 'Attack'), and a reward function that provides feedback based on classification correctness.
3. **Architectural Implementation:** Two Q-network architectures will be implemented in Python using deep learning libraries. The first will be a standard Multi-Layer Perceptron. The second will be a Convolutional Neural Network, which requires reshaping the one-dimensional feature vector into a two-dimensional "feature-image" to exploit spatial correlations between features.
4. **Training and Evaluation:** Both models will be trained using a novel loop that blends supervised signals (the true labels) with RL principles (value estimation of future states). The trained models will be evaluated on a held-out test set to assess their detection performance and generalization capabilities.

Thesis Structure

This thesis is organized into six chapters, each building upon the last to present a complete account of the research:

- **Chapter 1: Intrusion Detection Systems (IDS)** lays the theoretical groundwork by discussing fundamental security concepts, detailing the various types of DoS/DDoS attacks, and explaining the architecture and limitations of traditional IDS.
- **Chapter 2: Reinforcement Learning for Cybersecurity** introduces the core concepts of Reinforcement Learning, with a specific focus on Deep Q-Learning. It explores how the RL paradigm can be modeled for security tasks, covering state, action, and reward design.

- **Chapter 3: Architectural Blueprint and Implementation** presents the technical core of this work. It details the data preprocessing pipeline and the architectural design of the two proposed DQN agents (MLP and CNN-based).
- **Chapter 4: Experiments and Results** will describe the experimental setup, present the performance results of the trained models, and provide a comparative analysis and discussion of the findings.

Chapter 1

Intrusion Detection System (IDS)

1.1 Security Fundamentals

Introduction

Security fundamentals refer to the essential principles, concepts, and practices that form the foundation of information security. These fundamentals encompass a wide range of technical and organizational measures aimed at protecting sensitive information and systems from unauthorized access, theft, damage, or other forms of compromise.

Key security fundamentals include Confidentiality, Integrity, Availability, Authentication, Authorization, Encryption, Risk Management, Incident Response, and Disaster Recovery. Together, these principles establish the basis of a comprehensive information security program, enabling organizations to effectively safeguard their critical information assets and maintain the trust of stakeholders.[1]

1.1.1 Core Security Concepts

Network security encompasses strategies and technologies to protect systems from cyber threats, particularly Denial of Service (DoS) attacks that aim to disrupt service availability. Several core principles serve as the foundation for secure systems.

CIA Triad

Confidentiality Confidentiality ensures that sensitive information is only accessible to authorized users. Techniques such as encryption, user authentication, and access control policies prevent unauthorized data access. While DoS attacks do not typically aim to breach confidentiality directly, successful exploitation may lead to indirect confidentiality violations if attackers cause service misconfigurations or force failovers to insecure states.

Integrity Integrity guarantees that data is accurate and unaltered. It is maintained through cryptographic hash functions, checksums, and digital signatures that validate whether data has been tampered with. During a DoS attack, integrity may be compromised by interrupting legitimate updates or corrupting processes due to system overload.

Availability Availability ensures that services and systems remain accessible to legitimate users at all times. This principle is the primary target of DoS attacks, which flood networks or services with excessive requests, causing slowdowns or complete denial of access. Maintaining availability involves redundancy, load balancing, and proactive mitigation strategies like rate limiting and firewalls.



Figure 1.1: Principles of information security

Source: [2]

Extended Security Principles

Non-repudiation Non-repudiation guarantees that an entity cannot deny having performed a particular action, such as sending a message or initiating a transaction. This is enforced through digital signatures and secure logging mechanisms. In the context of DoS attacks, non-repudiation helps trace attack origins and supports legal accountability.

Authenticity Authenticity confirms that data, communications, or users are genuine and not forged. Authentication protocols, digital certificates, and cryptographic techniques are used to ensure that data comes from trusted sources. This is crucial for filtering legitimate traffic from spoofed attack traffic in DoS scenarios.

Accountability Accountability ensures that all actions within a system can be traced to responsible users or processes. It involves logging, auditing, and monitoring to track behavior. Accountability is vital for forensic analysis after a DoS attack and for strengthening defenses against future intrusions.

1.2 Denial-of-Service (DoS) Attacks

Introduction Denial-of-Service (DoS) attacks aim to make a system or network resource unavailable to its intended users by overwhelming it with excessive traffic or exploiting protocol-level vulnerabilities. These attacks can disrupt services, degrade performance, or completely shut down access. The most common types include:

1.2.1 Volumetric Attacks

Volumetric attacks aim to saturate a target's bandwidth by generating an overwhelming amount of traffic. This often involves amplification techniques or high-rate packet floods. Common examples include UDP floods and DNS amplification attacks, which leverage misconfigured servers to multiply traffic directed at the victim [3].

In a UDP flood, attackers send large numbers of spoofed UDP packets to random or specific ports, consuming bandwidth and processing power. DNS amplification uses small queries to open resolvers with the victim's IP address, causing them to return large responses to the victim.

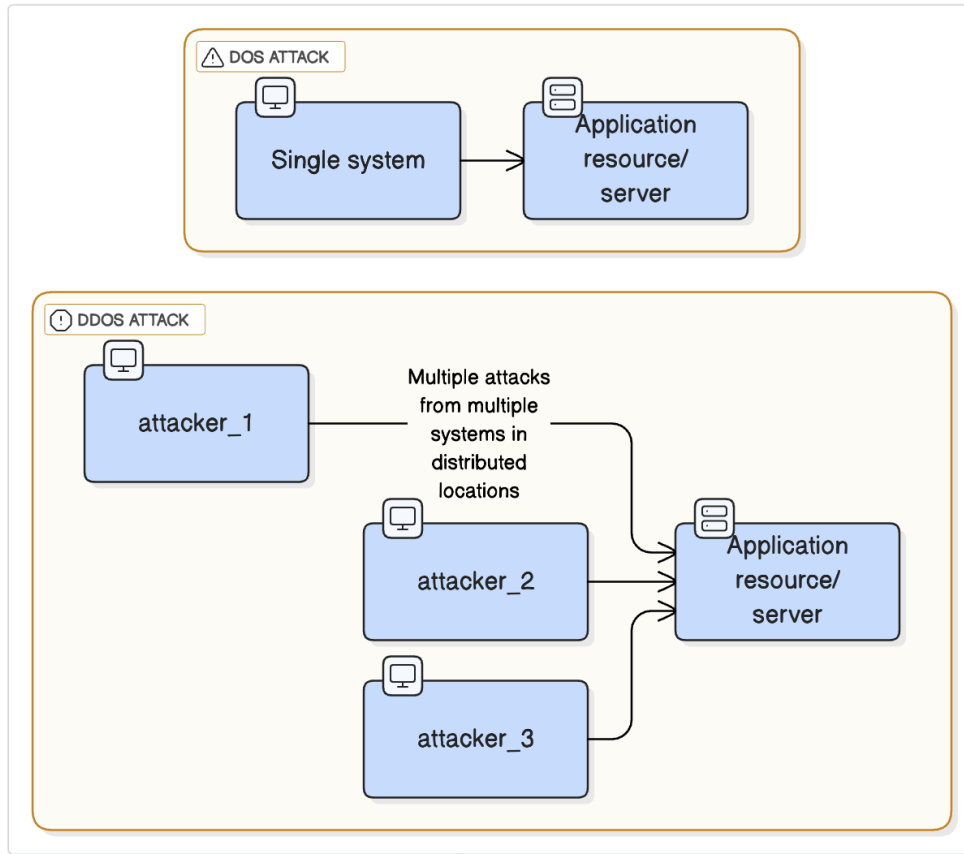


Figure 1.2: DoS and DDoS Attack Traffic Comparison

A notable example is the Mirai botnet, which infected hundreds of thousands of insecure IoT devices to coordinate massive traffic streams toward its targets [3]. The goal of volumetric attacks is to clog network links (Gbps or Tbps scale), preventing legitimate access.

1.2.2 Protocol Attacks

Protocol attacks exploit weaknesses in Layer 3/4 protocols (network or transport layer) to exhaust server or network resources. Unlike volumetric attacks, they do not require high bandwidth but instead consume stateful resources like connection tables or CPU cycles [3].

A classic example is the TCP SYN flood, where attackers send numerous SYN packets without completing the handshake, filling the server's connection queue [4]. ICMP floods and Ping of Death attacks exploit the Internet Control Message Protocol to crash or freeze systems by sending malformed or excessive traffic.

Other examples include fragmentation attacks (e.g., Teardrop), which send overlapping IP fragments to crash reassembly logic, and ACK/FIN floods that exhaust firewall state tables.

1.2.3 Application Layer Attacks

Application-layer (Layer 7) attacks generate traffic that appears legitimate at the application level but consumes excessive server resources such as CPU, memory, or threads [3]. These attacks are difficult to detect because they mimic real user behavior.

HTTP GET/POST floods can overwhelm web servers by triggering costly operations. A well-known example is Slowloris, which sends partial HTTP headers slowly to hold open many connections and exhaust the web server's pool [5]. Other examples include RUDY (R-U-Dead-Yet) attacks and HTTP floods targeting dynamic or database-backed content.

1.2.4 Distributed Denial-of-Service (DDoS)

DDoS attacks use multiple compromised machines (botnets) to launch coordinated attacks, making them harder to block and vastly more powerful than single-source DoS attacks [3].

Botnets like Mirai leverage IoT devices to simultaneously launch volumetric or protocol-based attacks. Reflective amplification (e.g., using open DNS/NTP servers) further increases traffic impact. Mitigation requires upstream filtering, rate-limiting, and often third-party scrubbing services.

1.2.5 Network Traffic Analysis

Mitigating DoS/DDoS attacks relies on monitoring and anomaly detection. This involves establishing a baseline of normal network behavior and identifying deviations in volume, protocol use, or source IPs [3, 6].

Anomaly-based intrusion detection systems (IDS) can flag unusual spikes, connection patterns, or TCP flag anomalies. Signature-based systems match known attack patterns. Flow analysis tools (e.g., NetFlow, sFlow) help detect irregular byte/packet rates or unusual source-destination pairs.

Advanced methods include machine learning to classify deviations. Once detected, defenses such as rate limiting, traffic shaping, or redirection to mitigation services are employed to maintain availability.

1.3 Intrusion Detection Systems (IDS)

An **Intrusion Detection System (IDS)** is a cybersecurity solution designed to monitor network or system activities for malicious actions or policy violations. Upon detecting such activities, the IDS typically alerts system administrators or integrates with centralized security tools like Security Information and Event Management (SIEM) systems to facilitate a coordinated response [7].

1.3.1 IDS Architecture

Monitored System Feedback Flow

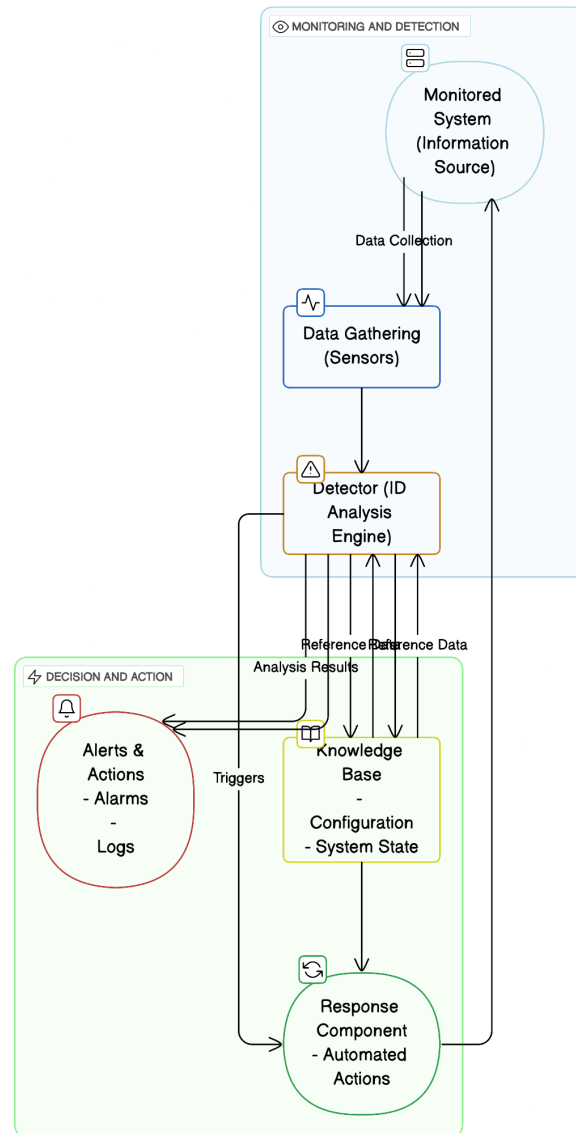


Figure 1.3: Network-based Intrusion Detection System (NIDS) Architecture

Modern IDS platforms implement modular architectures with several key components:

- **Sensors:** are the front line of an IDS, deployed at strategic network chokepoints—such as mirror (SPAN) ports, network taps, or virtual interfaces in cloud environments—to collect raw traffic data.
- **Analysis Engines:** form the IDS’s ”brain,” taking the sensor-collected data and applying a mixture of rule-based and intelligent techniques to detect threats.
- **Knowledge Base:** underpins all detection logic by maintaining up-to-date signatures, behavioral profiles, and historical datasets.
- **Response Systems:** close the loop between detection and defense. Once the analysis engine assigns a confidence score to an event, the response component generates alerts.

For DOS attack detection, these components must operate with high efficiency and scale to process enormous traffic volumes during attack scenarios.

1.3.2 IDS Types and Functionality

Intrusion Detection Systems serve as the primary monitoring and detection layer for network security threats, including DOS attacks:

- **Network-based IDS (NIDS):** A NIDS is a monitoring solution placed at key junctions within a network to passively capture and inspect all passing traffic.

NIDS (Network-based Intrusion Detection System)

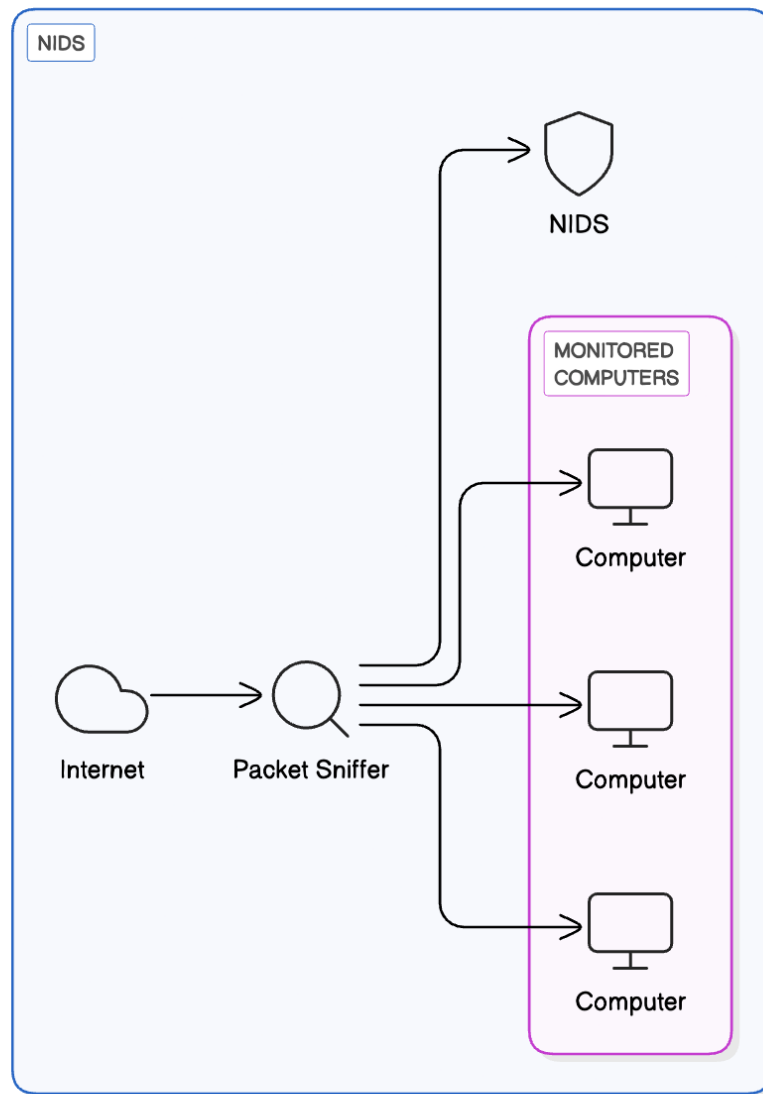


Figure 1.4: Network-based Intrusion Detection System (NIDS) Diagram

- **Host-based IDS (HIDS):** A HIDS operates at the level of individual endpoints or servers, offering deep visibility into system-specific activities.

HIDS Architecture Comparison

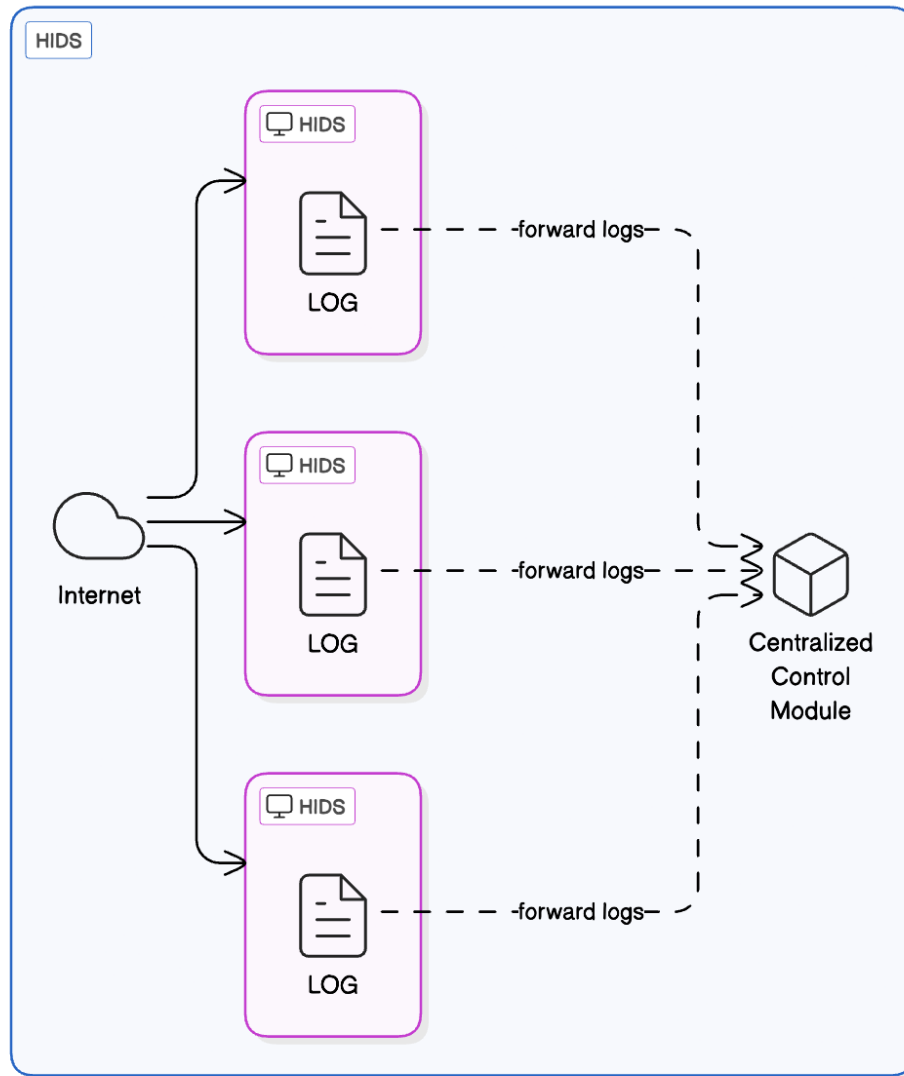


Figure 1.5: Host-based Intrusion Detection System (HIDS) Diagram

- **Detection Methodologies:**
 - **Signature-based:** Matches observed traffic against known attack patterns (signatures). Highly effective for known attacks but ineffective against zero-day attacks.
 - **Anomaly-based:** Establishes baselines of normal behavior and flags significant deviations. Can detect unknown attacks but may generate more false positives.
 - **Hybrid Systems:** Combine both methodologies to leverage their respective strengths and reduce false positives.

Limitations of Traditional IDS for DOS attack detection include difficulty processing high-volume traffic, distinguishing flash crowds from attacks, limited adaptation, high false-positive rates, and high resource consumption. These limitations necessitate AI-driven approaches like reinforcement learning.

1.4 Classification in Security Contexts

1.4.1 Classification Paradigms

Classification forms the foundation of automated threat detection: **Binary Classification:** Distinguishes between normal and attack traffic. **Multi-class Classification:** Distinguishes between multiple attack types. **Hierarchical Classification:** Structures threats in a tree-like format. **Multi-label Classification:** Assigns multiple categories to a single traffic flow.

What is Machine Learning?

Machine learning (ML) is a branch of artificial intelligence (AI) focused on enabling computers and machines to imitate the way that humans learn. The learning system consists of a decision process, an error function, and a model optimization process.

1.4.2 Machine Learning Techniques

Supervised Learning

Supervised learning involves training on labeled datasets. Common algorithms include Random Forests, Support Vector Machines (SVM), and Neural Networks.

Unsupervised Learning

Unsupervised learning focuses on uncovering patterns in unlabeled data. Techniques include clustering (k-means) and anomaly detection (Isolation Forest).

Semi-Supervised Learning

Semi-supervised learning combines a small amount of labeled data with a larger pool of unlabeled data, reducing the need for extensive manual labeling.

Chapter 2

Reinforcement Learning Applications: A Comprehensive Guide

2.1 Introduction to the Reinforcement Learning Approach

Reinforcement Learning (RL) is a branch of machine learning that focuses on how agents can learn to make decisions through trial and error to maximize cumulative rewards. Unlike supervised learning, where models learn from labeled data, RL involves learning optimal behaviors through interactions with an environment, receiving feedback in the form of rewards or penalties based on actions taken. This approach enables agents to discover strategies that yield the highest long-term benefits, making RL particularly effective for tasks involving sequential decision-making, such as robotics, game playing, and autonomous systems. Reinforcement Learning provides several key advantages over traditional DoS detection approaches:

Source: Adapted from GeeksforGeeks: What is Reinforcement Learning

2.2 Introduction to Deep Q-Learning

Deep Q-Learning is a powerful extension of the traditional Q-Learning algorithm that leverages deep neural networks to approximate the Q-value function. In standard Q-Learning, an agent maintains a Q-table that maps state-action pairs to expected future rewards. However, this becomes impractical for environments with large or continuous state spaces. Deep Q-Learning addresses this limitation by using a deep neural network, known as a Q-network, to estimate Q-values directly from raw input states.

In reinforcement learning, the agent interacts with the environment in discrete time steps. At each step t , the agent observes a state s_t , selects an action a_t , receives a reward r_t , and transitions to a new state s_{t+1} . The goal is to learn a policy π that maximizes the expected cumulative reward over time.

The Q-value function $Q(s, a)$ represents the expected return of taking action a in state s and following the policy thereafter. In Deep Q-Learning, the Q-network is trained to minimize the difference between the predicted Q-value and the target Q-value, which is computed using the Bellman equation:

$$Q(s_t, a_t) = r_t + \gamma \max_{a'} Q(s_{t+1}, a'; \theta^-) \quad (2.1)$$

Here, γ is the discount factor, θ are the parameters of the current Q-network, and θ^- are the parameters of a target network that is periodically updated to stabilize training.

To improve learning stability and efficiency, Deep Q-Learning introduces two key techniques:

- **Experience Replay:** Stores past experiences in a replay buffer and samples mini-batches randomly during training to break correlations between consecutive samples.
- **Target Network:** Uses a separate, periodically updated target network to compute the target Q-values, reducing oscillations and divergence.

Deep Q-Learning has been successfully applied to various complex decision-making tasks, such as playing Atari games directly from raw pixels, demonstrating its potential in handling high-dimensional input spaces.

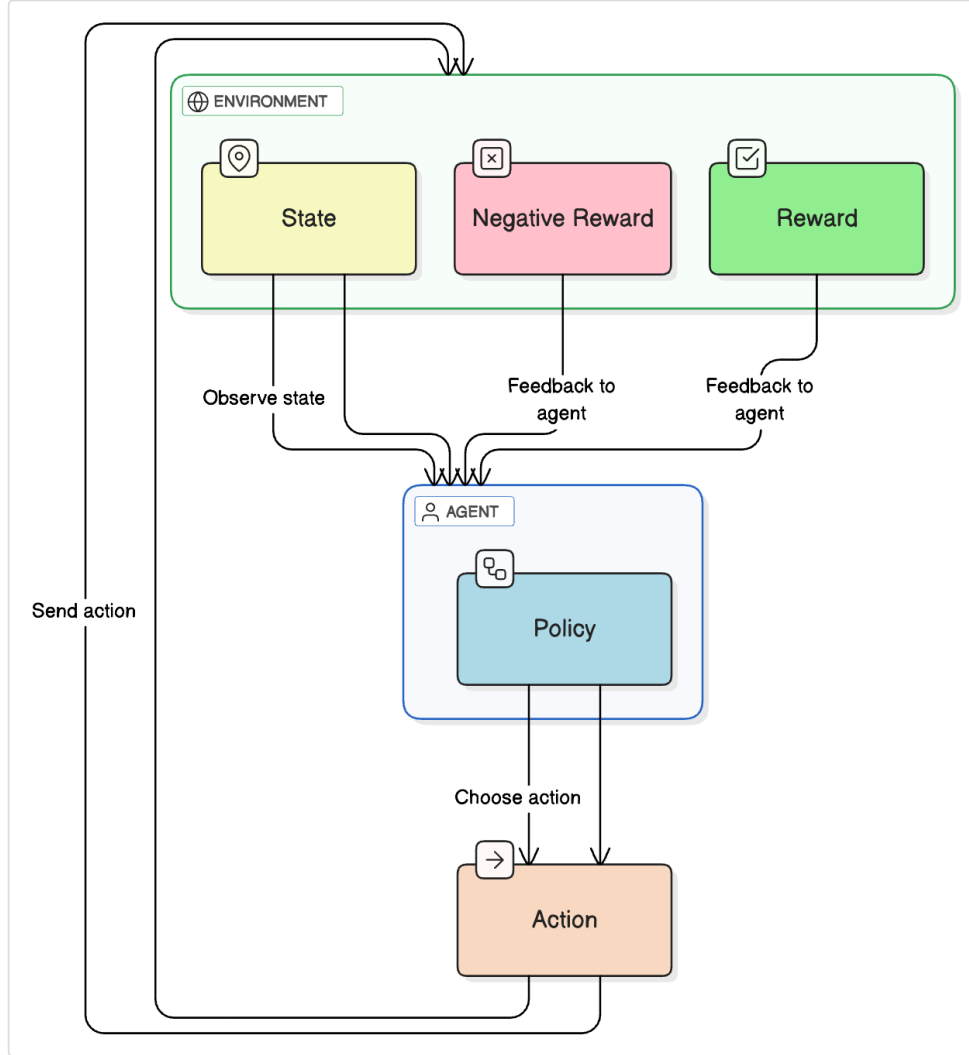


Figure: Basic Components of Reinforcement Learning

2.3 Environment Modeling for DoS Detection

The foundation of any RL application is a well-designed environment that accurately represents the problem domain. For DoS detection, this environment must model network traffic patterns and attack scenarios in a way that enables effective learning.

2.3.1 State Space Design

The state space in a reinforcement learning (RL) framework for Denial-of-Service (DoS) detection encodes the essential characteristics of network traffic that the agent uses to make decisions. A well-designed state space should reflect a comprehensive yet efficient representation of the network environment, enabling the RL agent to distinguish between normal and malicious activity. This design typically integrates four categories of features: traffic volume metrics, statistical distribution properties, temporal behavior indicators, and content-based descriptors. Each category offers a unique perspective on the nature of the traffic, contributing to a robust and informative state representation.

Traffic volume features quantify the scale and frequency of network activity. These include metrics such as packets per second, bytes per second, flow initiation rates, and the number of concurrent connections. These indicators are particularly useful for detecting volumetric DoS attacks, which are characterized by sudden spikes in traffic volume designed to overwhelm network resources. By monitoring these patterns, the RL agent can quickly recognize and react to abnormal surges in activity.

Statistical distribution features capture how traffic attributes are spread across various dimensions. Entropy values of source and destination IP addresses, protocol usage percentages, port distributions, and packet size variation are key indicators in this category. These features help identify anomalies that are not necessarily reflected in overall traffic volume but manifest as irregularities in distribution patterns, such as a high concentration of traffic targeting a specific port or originating from a narrow set of addresses.

Temporal pattern features provide insight into the evolution of traffic over time. By analyzing short-term trends, comparing current behavior to historical baselines, or evaluating periodicity, the agent can detect subtle or persistent attack patterns. Features such as time-of-day normalization or trend analysis allow the RL agent to differentiate between expected daily fluctuations and genuinely suspicious behavior, such as traffic bursts that align with known attack schedules.

Content-based features delve into the specifics of packet content and header information. This includes identifying protocol anomalies, unusual header field values, payload characteristics (when available), and application-layer request patterns. Such features are crucial for detecting more sophisticated attacks that may not exhibit volume or distribution anomalies but instead exploit specific protocol vulnerabilities or payload structures.

A critical aspect of state space design is balancing informativeness with efficiency. High-dimensional state representations can slow down learning and lead to overfitting. To mitigate this, techniques like feature selection, dimensionality reduction, and hierarchical representations can be employed. These methods help retain the most relevant features while reducing computational overhead, ultimately enhancing the agent’s learning performance and scalability.

2.3.2 Action Space Definition

The action space defines the set of all possible decisions or interventions that the reinforcement learning (RL) agent can make in response to observed network conditions. In the context of DoS detection systems, these actions are typically organized into three main categories: detection-related, response-related, and adaptive actions. Each category serves a distinct role in enabling the agent to not only identify malicious activity but also to react appropriately and adjust its behavior over time.

Detection-related actions allow the agent to classify the nature of incoming traffic. These may include flagging traffic as normal, marking it as a potential DoS attack with low confidence, confirming it as a high-confidence DoS event, or requesting additional inspection to reduce uncertainty. These actions contribute to the agent’s ability to differentiate between benign and malicious behavior with varying degrees of certainty.

Response-related actions are used to enforce protective measures against suspicious or confirmed attack traffic. These actions include allowing traffic to proceed unimpeded, rate-limiting suspected malicious flows, blocking specific traffic patterns known to be harmful, or redirecting traffic for more detailed analysis by external systems. Such actions form the core of the system’s active defense mechanism, enabling real-time mitigation of ongoing threats.

Adaptive actions provide the agent with the ability to adjust its operational parameters in response to evolving threat landscapes. This includes actions such as modifying monitoring sensitivity, tuning feature extraction processes, escalating suspicious cases to human analysts, or gathering additional contextual data to refine future decision-making. These adaptive strategies are essential for maintaining long-term performance in dynamic environments.

The granularity of the action space plays a critical role in determining both the learning efficiency and the practical utility of the RL system. An overly coarse action space may restrict the agent’s responsiveness and lead to underfitting, whereas an excessively fine-grained space can increase the complexity of the learning problem and slow convergence. Therefore, careful design of the action space is essential to strike an optimal balance between expressiveness and tractability.

2.3.3 Reward Function Design

The reward function plays a pivotal role in shaping the learning behavior of the reinforcement learning (RL) agent, guiding it toward effective and efficient DoS detection strategies. A well-crafted reward function must strike a balance between multiple, often competing objectives such as detection accuracy, operational efficiency, and timely response.

In terms of detection accuracy, the agent is incentivized through positive rewards for correctly identifying attacks, while being penalized for false positives and false negatives. These penalties and rewards may be further weighted based on the severity of the detected attack and the agent’s confidence in its classification. This ensures that the agent not only learns to detect attacks but also to assess their impact and act accordingly.

Operational efficiency is also integral to the reward function. Small penalties may be applied to account for resource consumption, including computational overhead and bandwidth usage. Additionally, time-based rewards can encourage early detection, while excessive inspection of normal traffic or inefficient use of mitigation mechanisms can incur penalties. These components help align the agent’s behavior with real-world constraints and performance expectations.

A typical structure for the reward function may take the following form:

$$R(s, a, s') = w_1 \times \text{DetectionAccuracy} + w_2 \times \text{TimeToDetect} + w_3 \times \text{ResourceEfficiency} \quad (2.2)$$

In this formulation, *DetectionAccuracy* reflects the correctness of the classification, considering true positives, false positives, and false negatives. *TimeToDetect* quantifies the latency in identifying an attack, and *ResourceEfficiency* evaluates how judiciously the system utilizes its resources during detection and response. The weights w_1, w_2, w_3 are tunable parameters that allow the system to prioritize objectives based on organizational or operational goals.

Careful calibration of the reward function is essential to ensure that the agent develops policies that are not only accurate and robust but also practical within the resource constraints and real-time requirements of modern network environments.

2.3.4 Environment Dynamics

The environment dynamics define the mechanisms through which state transitions occur in response to the agent’s actions. In the context of DoS detection, these dynamics are influenced by both the underlying traffic patterns and the mitigation strategies employed by the agent. State transitions are often governed by probabilistic rules that reflect how network traffic evolves over time. For instance, specific traffic behaviors may be more or less likely to follow certain patterns, such as bursty traffic or sustained high-volume flows indicative of an ongoing attack.

The agent’s actions directly impact the state evolution. For example, blocking suspicious traffic may lead to a reduction in attack volume, altering the observable state characteristics. Similarly, redirecting traffic for analysis or applying rate limits can change flow distributions and entropy metrics within the network. It is also important to consider temporal aspects of the environment; the effects of certain actions may not be immediately visible. For instance, a rate-limiting policy might only manifest noticeable changes after a delay, as the network adapts or the attacker responds.

Modeling these dynamics accurately is essential for realistic simulation and effective training of reinforcement learning agents. It allows the agent to learn how its actions influence the environment and adjust its policy accordingly to achieve robust and adaptive DoS detection performance.

2.3.5 Challenges in Applying RL to Security Domains

While reinforcement learning (RL) holds significant potential for enhancing cybersecurity systems, its application to security contexts—particularly DoS detection—presents several unique challenges. One major hurdle is the high dimensionality of the state space. Network traffic generates vast and complex feature sets, requiring sophisticated representation and dimensionality reduction techniques to ensure tractable learning.

Another difficulty lies in the delayed nature of rewards. In many cases, the outcome of a security-related decision—whether it successfully prevented an attack or caused unintended side effects—may only

become clear after a considerable time delay. This complicates credit assignment and policy evaluation. Moreover, rewards in the security domain tend to be sparse; attack events are relatively rare in comparison to the continuous stream of benign network activity, making it difficult for the agent to gain useful feedback consistently.

Safety during the learning process is also a critical concern. Exploration, a core part of RL, must be managed carefully to avoid compromising system integrity. Unlike in traditional domains, erroneous actions in security systems can have severe consequences. Compounding this is the adversarial nature of the environment—attackers actively adapt their strategies to bypass detection mechanisms, forcing RL agents to learn in a moving target setting.

To address these challenges, structured educational frameworks and simulation environments can be employed. These offer controlled scenarios where RL techniques can be gradually introduced and refined before deployment in production systems.

2.3.6 Advantages of RL for DoS Detection

Despite the complexities involved, RL brings numerous advantages to the domain of DoS detection. One of the most notable is adaptability. Unlike rule-based systems that require manual updates, RL agents can dynamically adjust to novel attack strategies and traffic behaviors. This makes them particularly well-suited for rapidly evolving threat landscapes.

RL also enables contextual decision-making, where responses are informed by the broader network state rather than isolated events. This holistic approach helps reduce false positives and improve detection accuracy. Moreover, RL frameworks can be designed to balance multiple, sometimes conflicting objectives—such as maximizing detection rates while minimizing the impact on normal traffic—through multi-objective optimization.

A proactive defense posture is another key benefit. By interacting continuously with the environment, RL agents can learn to anticipate and preempt attack progression rather than merely reacting to already-occurred incidents. This ongoing interaction fosters continuous improvement, allowing the system to refine its policies over time and maintain effectiveness even as network conditions and attack techniques evolve.

2.4 Introduction to the Reinforcement Learning Approach

Reinforcement Learning (RL) is a branch of machine learning that focuses on how agents can learn to make decisions through trial and error to maximize cumulative rewards. Unlike supervised learning, where models learn from labeled data, RL involves learning optimal behaviors through interactions with an environment, receiving feedback in the form of rewards or penalties based on actions taken. This approach enables agents to discover strategies that yield the highest long-term benefits, making RL particularly effective for tasks involving sequential decision-making, such as robotics, game playing, and autonomous systems. Reinforcement Learning provides several key advantages over traditional DoS detection approaches:

Source: Adapted from GeeksforGeeks: What is Reinforcement Learning

2.5 Introduction to Deep Q-Learning

Deep Q-Learning is a powerful extension of the traditional Q-Learning algorithm that leverages deep neural networks to approximate the Q-value function. In standard Q-Learning, an agent maintains a Q-table that maps state-action pairs to expected future rewards. However, this becomes impractical for environments with large or continuous state spaces. Deep Q-Learning addresses this limitation by using a deep neural network, known as a Q-network, to estimate Q-values directly from raw input states.

In reinforcement learning, the agent interacts with the environment in discrete time steps. At each step t , the agent observes a state s_t , selects an action a_t , receives a reward r_t , and transitions to a new state s_{t+1} . The goal is to learn a policy π that maximizes the expected cumulative reward over time.

The Q-value function $Q(s, a)$ represents the expected return of taking action a in state s and following the policy thereafter. In Deep Q-Learning, the Q-network is trained to minimize the difference between the predicted Q-value and the target Q-value, which is computed using the Bellman equation:

$$Q(s_t, a_t) = r_t + \gamma \max_{a'} Q(s_{t+1}, a'; \theta^-) \quad (2.3)$$

Here, γ is the discount factor, θ are the parameters of the current Q-network, and θ^- are the parameters of a target network that is periodically updated to stabilize training.

To improve learning stability and efficiency, Deep Q-Learning introduces two key techniques:

- **Experience Replay:** Stores past experiences in a replay buffer and samples mini-batches randomly during training to break correlations between consecutive samples.
- **Target Network:** Uses a separate, periodically updated target network to compute the target Q-values, reducing oscillations and divergence.

Deep Q-Learning has been successfully applied to various complex decision-making tasks, such as playing Atari games directly from raw pixels, demonstrating its potential in handling high-dimensional input spaces.

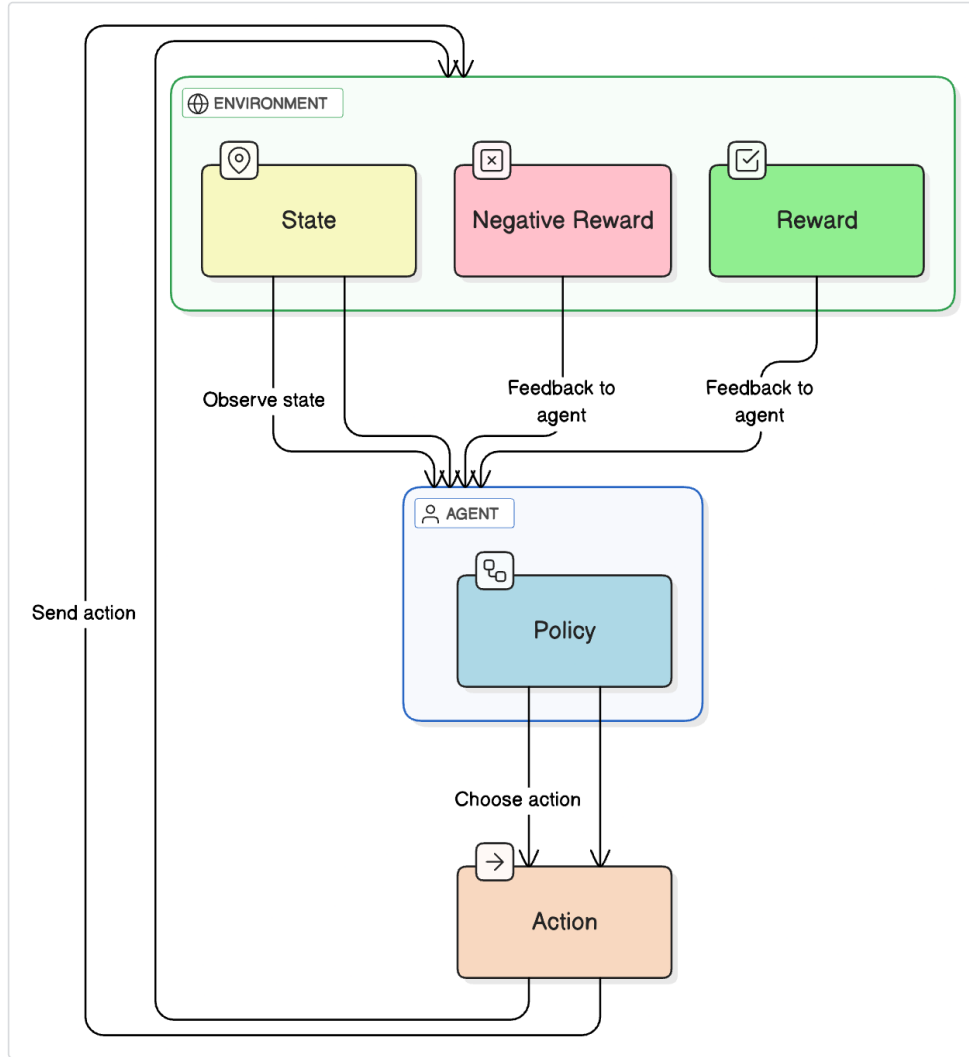


Figure: Basic Components of Reinforcement Learning

2.6 Environment Modeling for DoS Detection

The foundation of any RL application is a well-designed environment that accurately represents the problem domain. For DoS detection, this environment must model network traffic patterns and attack scenarios in a

way that enables effective learning.

2.6.1 State Space Design

The state space in a reinforcement learning (RL) framework for Denial-of-Service (DoS) detection encodes the essential characteristics of network traffic that the agent uses to make decisions. A well-designed state space should reflect a comprehensive yet efficient representation of the network environment, enabling the RL agent to distinguish between normal and malicious activity. This design typically integrates four categories of features: traffic volume metrics, statistical distribution properties, temporal behavior indicators, and content-based descriptors. Each category offers a unique perspective on the nature of the traffic, contributing to a robust and informative state representation.

Traffic volume features quantify the scale and frequency of network activity. These include metrics such as packets per second, bytes per second, flow initiation rates, and the number of concurrent connections. These indicators are particularly useful for detecting volumetric DoS attacks, which are characterized by sudden spikes in traffic volume designed to overwhelm network resources. By monitoring these patterns, the RL agent can quickly recognize and react to abnormal surges in activity.

Statistical distribution features capture how traffic attributes are spread across various dimensions. Entropy values of source and destination IP addresses, protocol usage percentages, port distributions, and packet size variation are key indicators in this category. These features help identify anomalies that are not necessarily reflected in overall traffic volume but manifest as irregularities in distribution patterns, such as a high concentration of traffic targeting a specific port or originating from a narrow set of addresses.

Temporal pattern features provide insight into the evolution of traffic over time. By analyzing short-term trends, comparing current behavior to historical baselines, or evaluating periodicity, the agent can detect subtle or persistent attack patterns. Features such as time-of-day normalization or trend analysis allow the RL agent to differentiate between expected daily fluctuations and genuinely suspicious behavior, such as traffic bursts that align with known attack schedules.

Content-based features delve into the specifics of packet content and header information. This includes identifying protocol anomalies, unusual header field values, payload characteristics (when available), and application-layer request patterns. Such features are crucial for detecting more sophisticated attacks that may not exhibit volume or distribution anomalies but instead exploit specific protocol vulnerabilities or payload structures.

A critical aspect of state space design is balancing informativeness with efficiency. High-dimensional state representations can slow down learning and lead to overfitting. To mitigate this, techniques like feature selection, dimensionality reduction, and hierarchical representations can be employed. These methods help retain the most relevant features while reducing computational overhead, ultimately enhancing the agent's learning performance and scalability.

2.6.2 Action Space Definition

The action space defines the set of all possible decisions or interventions that the reinforcement learning (RL) agent can make in response to observed network conditions. In the context of DoS detection systems, these actions are typically organized into three main categories: detection-related, response-related, and adaptive actions. Each category serves a distinct role in enabling the agent to not only identify malicious activity but also to react appropriately and adjust its behavior over time.

Detection-related actions allow the agent to classify the nature of incoming traffic. These may include flagging traffic as normal, marking it as a potential DoS attack with low confidence, confirming it as a high-confidence DoS event, or requesting additional inspection to reduce uncertainty. These actions contribute to the agent's ability to differentiate between benign and malicious behavior with varying degrees of certainty.

Response-related actions are used to enforce protective measures against suspicious or confirmed attack traffic. These actions include allowing traffic to proceed unimpeded, rate-limiting suspected malicious flows, blocking specific traffic patterns known to be harmful, or redirecting traffic for more detailed analysis by external systems. Such actions form the core of the system's active defense mechanism, enabling real-time mitigation of ongoing threats.

Adaptive actions provide the agent with the ability to adjust its operational parameters in response to evolving threat landscapes. This includes actions such as modifying monitoring sensitivity, tuning feature extraction processes, escalating suspicious cases to human analysts, or gathering additional contextual data to refine future decision-making. These adaptive strategies are essential for maintaining long-term performance in dynamic environments.

The granularity of the action space plays a critical role in determining both the learning efficiency and the practical utility of the RL system. An overly coarse action space may restrict the agent’s responsiveness and lead to underfitting, whereas an excessively fine-grained space can increase the complexity of the learning problem and slow convergence. Therefore, careful design of the action space is essential to strike an optimal balance between expressiveness and tractability.

2.6.3 Reward Function Design

The reward function plays a pivotal role in shaping the learning behavior of the reinforcement learning (RL) agent, guiding it toward effective and efficient DoS detection strategies. A well-crafted reward function must strike a balance between multiple, often competing objectives such as detection accuracy, operational efficiency, and timely response.

In terms of detection accuracy, the agent is incentivized through positive rewards for correctly identifying attacks, while being penalized for false positives and false negatives. These penalties and rewards may be further weighted based on the severity of the detected attack and the agent’s confidence in its classification. This ensures that the agent not only learns to detect attacks but also to assess their impact and act accordingly.

Operational efficiency is also integral to the reward function. Small penalties may be applied to account for resource consumption, including computational overhead and bandwidth usage. Additionally, time-based rewards can encourage early detection, while excessive inspection of normal traffic or inefficient use of mitigation mechanisms can incur penalties. These components help align the agent’s behavior with real-world constraints and performance expectations.

A typical structure for the reward function may take the following form:

$$R(s, a, s') = w_1 \times \text{DetectionAccuracy} + w_2 \times \text{TimeToDetect} + w_3 \times \text{ResourceEfficiency} \quad (2.4)$$

In this formulation, *DetectionAccuracy* reflects the correctness of the classification, considering true positives, false positives, and false negatives. *TimeToDetect* quantifies the latency in identifying an attack, and *ResourceEfficiency* evaluates how judiciously the system utilizes its resources during detection and response. The weights w_1, w_2, w_3 are tunable parameters that allow the system to prioritize objectives based on organizational or operational goals.

Careful calibration of the reward function is essential to ensure that the agent develops policies that are not only accurate and robust but also practical within the resource constraints and real-time requirements of modern network environments.

2.6.4 Environment Dynamics

The environment dynamics define the mechanisms through which state transitions occur in response to the agent’s actions. In the context of DoS detection, these dynamics are influenced by both the underlying traffic patterns and the mitigation strategies employed by the agent. State transitions are often governed by probabilistic rules that reflect how network traffic evolves over time. For instance, specific traffic behaviors may be more or less likely to follow certain patterns, such as bursty traffic or sustained high-volume flows indicative of an ongoing attack.

The agent’s actions directly impact the state evolution. For example, blocking suspicious traffic may lead to a reduction in attack volume, altering the observable state characteristics. Similarly, redirecting traffic for analysis or applying rate limits can change flow distributions and entropy metrics within the network. It is also important to consider temporal aspects of the environment; the effects of certain actions may not be immediately visible. For instance, a rate-limiting policy might only manifest noticeable changes after a delay, as the network adapts or the attacker responds.

Modeling these dynamics accurately is essential for realistic simulation and effective training of reinforcement learning agents. It allows the agent to learn how its actions influence the environment and adjust its policy accordingly to achieve robust and adaptive DoS detection performance.

2.6.5 Challenges in Applying RL to Security Domains

While reinforcement learning (RL) holds significant potential for enhancing cybersecurity systems, its application to security contexts—particularly DoS detection—presents several unique challenges. One major hurdle is the high dimensionality of the state space. Network traffic generates vast and complex feature sets, requiring sophisticated representation and dimensionality reduction techniques to ensure tractable learning.

Another difficulty lies in the delayed nature of rewards. In many cases, the outcome of a security-related decision—whether it successfully prevented an attack or caused unintended side effects—may only become clear after a considerable time delay. This complicates credit assignment and policy evaluation. Moreover, rewards in the security domain tend to be sparse; attack events are relatively rare in comparison to the continuous stream of benign network activity, making it difficult for the agent to gain useful feedback consistently.

Safety during the learning process is also a critical concern. Exploration, a core part of RL, must be managed carefully to avoid compromising system integrity. Unlike in traditional domains, erroneous actions in security systems can have severe consequences. Compounding this is the adversarial nature of the environment—attackers actively adapt their strategies to bypass detection mechanisms, forcing RL agents to learn in a moving target setting.

To address these challenges, structured educational frameworks and simulation environments can be employed. These offer controlled scenarios where RL techniques can be gradually introduced and refined before deployment in production systems.

2.6.6 Advantages of RL for DoS Detection

Despite the complexities involved, RL brings numerous advantages to the domain of DoS detection. One of the most notable is adaptability. Unlike rule-based systems that require manual updates, RL agents can dynamically adjust to novel attack strategies and traffic behaviors. This makes them particularly well-suited for rapidly evolving threat landscapes.

RL also enables contextual decision-making, where responses are informed by the broader network state rather than isolated events. This holistic approach helps reduce false positives and improve detection accuracy. Moreover, RL frameworks can be designed to balance multiple, sometimes conflicting objectives—such as maximizing detection rates while minimizing the impact on normal traffic—through multi-objective optimization.

A proactive defense posture is another key benefit. By interacting continuously with the environment, RL agents can learn to anticipate and preempt attack progression rather than merely reacting to already-occurred incidents. This ongoing interaction fosters continuous improvement, allowing the system to refine its policies over time and maintain effectiveness even as network conditions and attack techniques evolve.

Chapter 3

Designing DQN Architectures for DoS/DDoS Attack Detection

3.1 Introduction

This chapter details the design and conceptualization of a robust Distributed Denial of Service (DDoS) attack detection system. We propose an innovative approach that leverages the power of Deep Q-Networks (DQN), a prominent algorithm in Reinforcement Learning (RL), by adapting it to the context of a supervised classification problem. The primary goal is to develop an intelligent agent capable of accurately identifying DoS/DDoS attacks in real-world network traffic scenarios.

The chapter begins by meticulously defining the problem within the RL framework, mapping traditional RL concepts such as states, actions, and rewards to the specifics of network intrusion detection. We then outline the comprehensive data preprocessing pipeline, which is crucial for preparing raw network traffic data for effective model training. Subsequently, we delve into the proposed Deep Q-Network (DQN) architectures, including discussions on the standard DQN and the enhanced Double DQN (DDQN), as well as integrating Convolutional Neural Networks (CNN) for advanced feature extraction. The operational mechanics of the DQN agent, encompassing the training loop and its core components, are thoroughly explained. This foundation lays the groundwork for the practical implementation and evaluation of these architectures in subsequent chapters.

3.2 Problem Modeling for DoS/DDoS Detection using DQN

To leverage Reinforcement Learning (RL) for DDoS attack detection, we must reframe the classification task into an RL problem. This involves defining key RL components: the environment, rewards, agent's objective, and classification tasks within our context.

3.2.1 Environment Definition

In our framework, the "environment" is simulated by a preprocessed dataset. Each state (s_t) is a single, preprocessed feature vector representing a snapshot of network traffic. Unlike traditional RL where states evolve from prior actions, here, the "next state" (s_{t+1}) is functionally defined as the subsequent sample in a randomly shuffled batch. This approach allows us to estimate future value for Q-learning updates, even in a non-sequential data context. The dataset itself acts as the source of observations the agent makes at each "time step."

Reward Function

The reward (r_t) is an immediate scalar feedback signal provided by the simulated environment based on the agent's classification decision. It quantifies the quality of the agent's chosen action (a_t) by directly comparing

it with the known optimal action (a_t^*), i.e., the true label.

The reward function is defined as:

$$r_t = \begin{cases} +10 & \text{if } a_t = a_t^* \quad (\text{Correct Classification}) \\ 0 & \text{if } a_t \neq a_t^* \quad (\text{Incorrect Classification}) \end{cases}$$

Using a neutral reward of 0 for incorrect actions, rather than a negative penalty (e.g., -1), can lead to more stable training in some contexts, as it does not excessively punish the agent during its initial exploration phase. This reward mechanism forms the crucial bridge between the supervised ground-truth label and the reinforcement learning feedback loop.

Agent’s Objective

The primary objective of the reinforcement learning (RL) agent is to learn an optimal detection policy $\pi(a | s)$ that maximizes its long-term cumulative reward.

In the context of DDoS detection, this objective translates into maximizing the number of correct classifications—that is, accurately identifying both benign and malicious traffic.

By continuously refining its understanding of the value associated with selecting a particular label in a given context, the agent aims to converge toward a strategy that minimizes misclassifications and ensures robust detection performance.

Definition of Classification Tasks

The problem is framed as a binary classification task to simplify the detection process. The agent’s actions (a_t) correspond to the two possible classification labels:

- **Action 0:** Classify the network traffic as **Benign**.
- **Action 1:** Classify the network traffic as **Attack** (encompassing any type of DoS/DDoS attack).

The multi-class labels present in the original dataset (e.g., **DrDoS_NTP**, **SYN**, **UDP-lag**) are converted into this binary format during preprocessing. This transformation enables the model to focus on the fundamental task of distinguishing any malicious activity from normal network behavior.

3.2.2 Data Preprocessing Pipeline

The preprocessing phase is a fundamental component of building an effective machine learning model, ensuring that the input data is clean, consistent, and optimally formatted for the learning algorithm. Our approach involves several key steps:

Dataset

The primary dataset used for this research is a publicly available dataset. This dataset is highly relevant as it contains a variety of the most up-to-date common DDoS attacks alongside benign traffic, closely resembling real-world network scenarios. The data is provided in PCAP format and has been pre-processed using a specialized tool to generate labeled flows in CSV files. Each flow is characterized by features based on its time stamp, source and destination IPs, source and destination ports, and the protocol used.

Step 1: Data Cleaning

The initial phase focuses on cleaning the dataset to ensure its quality and integrity. This involves three key actions:

- **Handling Null Values:** The dataset is first inspected for any missing or null values. Rows containing such values are dropped entirely. This step is crucial to prevent computational errors during model training and to ensure that the model learns from complete, high-quality samples.

- **Dropping Irrelevant Features:** Certain features that are unique to individual flows and do not contribute to generalizable attack patterns are removed. Including them would introduce noise and create a high risk of overfitting, where the model memorizes specific instances from the training data rather than learning the underlying behavior of attacks.
- **Numeric Conversion:** All remaining feature values are converted to a numeric data type. Machine learning models, and particularly neural networks like the DQN, require numerical input for their mathematical operations. This step ensures that all data, including any encoded categorical features, is in a format that the model can process.

Step 2: Feature Selection using Random Forest Classifier

Feature selection is a critical process for reducing the dimensionality of the dataset, which in turn reduces model complexity, decreases training time, and mitigates the risk of overfitting. For this task, we utilize the *Random Forest Classifier*.

Random Forest is an ensemble learning method that constructs a multitude of decision trees during training. It provides a built-in measure of feature importance, which reflects the degree to which each feature contributes to improving the purity of the nodes in the trees. The process is as follows:

- A Random Forest Classifier is trained on the cleaned dataset.
- The `feature_importances_` attribute of the trained model is used to extract an importance score for each feature.
- Features are ranked based on these scores, and only the top-ranked features that contribute most significantly to the model’s predictive power are retained for subsequent steps.

Step 3: Data Normalization

After feature selection, the data is normalized to ensure that all features contribute equally to the model’s learning process. We employ *Min-Max normalization*, which scales each feature to a fixed range between 0 and 1.

This technique is essential for neural networks, as it prevents features with larger numeric ranges from dominating the learning process and helps stabilize the gradient descent optimization, leading to faster convergence.

The Min-Max normalization formula for a feature X is given by:

$$X_{\text{normalized}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

where:

- $X_{\text{normalized}}$ is the scaled value,
- X is the original feature value,
- X_{\min} is the minimum value of the feature in the dataset,
- X_{\max} is the maximum value of the feature in the dataset.

Step 4: Dataset Balancing with Random Undersampling

The dataset exhibits a significant class imbalance, with a much larger proportion of attack traffic compared to benign traffic. Training a model on such an imbalanced dataset would bias it towards the majority class (attacks), leading to poor detection performance for the minority class (benign traffic) and a high false positive rate.

To address this, we apply *random undersampling*. This technique balances the class distribution by reducing the number of samples from the majority class. Specifically, we randomly select and remove samples from the ‘attack’ class until its size matches the number of samples in the ‘benign’ class.

While this method risks discarding potentially useful information from the majority class, it is effective in creating a balanced dataset that enables the model to learn the distinguishing characteristics of both classes with equal importance. The result is a model that is more robust and accurate in classifying both attack and benign traffic.

3.2.3 Proposed DQN Architectures for Attack Detection

Having established the conceptual mapping between reinforcement learning and supervised classification, this section details the operational mechanics of the Deep Q-Network (DQN) agent and its proposed architectural implementations.

We begin by providing a high-level, component-wise breakdown of the unified training loop that forms the core of our methodology. This framework is *architecturally agnostic*, meaning it governs the learning process regardless of the specific neural network used for the internal Q-network.

The primary goal is to iteratively refine the agent’s Q-network, which is responsible for approximating the optimal action-value function $Q^*(s, a)$. This function estimates the long-term value of taking a specific action (i.e., making a classification) when presented with a particular state (i.e., a network traffic sample).

Deep Q-Learning Architecture (DQN)

The standard Deep Q-Learning (DQN) architecture is employed to learn the optimal Q-values. The training process involves continuous interaction with the simulated environment (the dataset) and iterative updates to the Q-network.

Component 1: Generic Sample and Data Provisioning

The learning process begins with the provisioning of data to the agent. At each training step, a generic sample is drawn from the preprocessed and shuffled dataset. This sample is a tuple containing three critical elements required for the Q-learning update: (s_t, a_t^*, s_{t+1}) .

- **Current State** (s_t): The feature vector of the network traffic sample currently under evaluation.
- **Optimal Action** (a_t^*): The ground-truth label corresponding to state s_t , re-conceptualized as the “optimal action.”
- **Next State** (s_{t+1}): The feature vector of the subsequent sample in the shuffled dataset, serving as a lookahead mechanism to estimate the value of the next state.

This data sampling process is formally represented by Algorithm 1.

Algorithm 1 Data Sampling Step

- 1: // Retrieve a minibatch of N generic samples from the pre-structured dataset.
 - 2: // Each sample is a tuple (s_t, a_t^*, s_{t+1}) .
 - 3: minibatch \leftarrow SampleFromDataset(dataset, batch_size = N)
 - 4: // Unpack the minibatch into separate tensors for processing.
 - 5: $S_t, A_t^*, S_{t+1} \leftarrow$ unpack(minibatch)
 - 6: // S_t : Current state vectors.
 - 7: // A_t^* : Ground-truth label indices (optimal actions).
 - 8: // S_{t+1} : Next state vectors.
-

Component 2: The ϵ -Greedy Policy and Action Selection

Once presented with the current state s_t , the agent must select an action a_t . This decision is governed by an ϵ -greedy policy—a fundamental strategy in reinforcement learning designed to manage the critical exploration-exploitation trade-off.

- **Exploitation:** With probability $(1 - \epsilon)$, the agent chooses to exploit its current knowledge. It feeds the state s_t into its Q-network, which predicts the Q-value for each possible action. The agent then greedily selects the action with the highest predicted Q-value:

$$a_t = \arg \max_a Q(s_t, a)$$

- **Exploration:** With probability ϵ , the agent explores. It disregards its learned policy and selects an action at random from the available action space.

The parameter ϵ is a tunable hyperparameter that is typically annealed over the course of training, gradually decaying from a high initial value to a smaller value. This process is formally represented by Algorithm 2.

Algorithm 2 Epsilon-Greedy Action Selection

```

1: function SELECTACTION(state,  $q\_network$ ,  $\epsilon$ )
2:    $p \leftarrow \text{random\_float}(0, 1)$ 
3:   if  $p < \epsilon$  then
4:     return random_choice(ACTION_SPACE)
5:   else
6:      $q\_values \leftarrow q\_network.\text{predict}(\text{state})$ 
7:     return  $\arg \max_a q\_values[a]$ 
8:   end if
9: end function

```

Component 3: The Target Q-Value (q_{ref}) Calculation

This step is one of the most sophisticated and innovative aspects of the entire training framework. In our DQN framework, the Q-network is trained against a dynamically constructed hybrid target vector, q_{ref} . This target ingeniously blends the explicit signal from the supervised label with the value-estimation principles of reinforcement learning.

The construction of this target vector for the current state s_t proceeds as follows:

- **Estimate Optimal Future Value:** The Q-network is used to predict the Q-values for the next state s_{t+1} . The maximum of these values is taken as the best estimate of the optimal value achievable from the next state:

$$\max_a Q(s_{t+1}, a)$$

- **Initialize Target Vector:** A target vector q_{ref} , of the same dimension as the action space, is initialized with zeros:

$$q_{\text{ref}} = \mathbf{0}$$

- **Inject the Supervised Signal:** A value of 1.0 is placed in the target vector at the index corresponding to the true label a_t^* :

$$q_{\text{ref}}[a_t^*] = 1.0$$

- **Inject the Reinforcement Learning Signal:** The estimated optimal future value is discounted by a factor λ and added to the index of the agent's chosen action a_t :

$$q_{\text{ref}}[a_t] += \lambda \cdot \max_a Q(s_{t+1}, a)$$

This dual-signal approach ensures the agent learns not only to identify the correct answer but also to accurately value its own decisions in the context of future possibilities. The entire process is formally described in Algorithm 3.

Algorithm 3 Target Q-Value Calculation

```
1: function CALCULATETARGET( $q\_network$ ,  $S_t$ ,  $S_{t+1}$ ,  $A^*$ ,  $A$ ,  $\lambda$ )
2:    $q\_next \leftarrow q\_network.predict(S_{t+1})$ 
3:    $max\_q\_next \leftarrow \max(q\_next, axis = 1)$ 
4:    $q\_ref \leftarrow \text{zeros\_like}(q\_network.predict(S_t))$ 
5:   for all  $i$  in batch do
6:      $true\_label \leftarrow A^*[i]$ 
7:      $agent\_action \leftarrow A[i]$ 
8:      $discounted\_future \leftarrow \lambda \cdot max\_q\_next[i]$ 
9:      $q\_ref[i][true\_label] \leftarrow 1.0$ 
10:     $q\_ref[i][agent\_action] \leftarrow q\_ref[i][agent\_action] + discounted\_future$ 
11:   end for
12:   return  $q\_ref$ 
13: end function
```

Component 4: Loss Calculation and Network Update

The final step in the training loop is to update the Q-network’s parameters (weights and biases) to improve its predictions. This is achieved by quantifying the error between the network’s output and the hybrid target, and then adjusting the weights to minimize this error.

- **Final Target Preparation:** The hybrid target vector q_{ref} may contain multiple non-zero values with continuous components. To make it compatible with a classification-style loss function, it is first converted into a definitive one-hot vector. This is done by finding the index of the maximum value in q_{ref} :

$$i^* = \arg \max_i q_{ref}[i]$$

The one-hot target vector \hat{y} is then constructed such that:

$$\hat{y}[i] = \begin{cases} 1 & \text{if } i = i^* \\ 0 & \text{otherwise} \end{cases}$$

- **Loss Function:** The discrepancy between the Q-network’s predicted Q-values (after applying the Softmax function to obtain a probability distribution) and the one-hot vector \hat{y} is computed using the Cross-Entropy Loss:

$$\mathcal{L} = - \sum_i \hat{y}[i] \cdot \log(\text{Softmax}(Q(s_t))[i])$$

- **Optimization:** The computed loss \mathcal{L} is used to guide the optimization process. An optimizer such as Adam is employed to compute the gradients of the loss with respect to the network’s parameters.
- **Backpropagation:** These gradients are propagated backward through the layers of the Q-network using the backpropagation algorithm, from the output layer toward the input.
- **Weight Update:** The optimizer then updates the parameters of the Q-network in the direction that minimizes the loss. This parameter update is denoted by:

$$Q^{(t)} \rightarrow Q^{(t+1)}$$

This complete update cycle forms a single iteration of the learning process. Repeating this process over many training samples and epochs enables the Q-network to gradually converge to an optimal policy, effectively learning to estimate the value of each classification choice for any given network traffic profile.

The architectural details of the Q-network will be discussed in the following subsections, specifically outlining both a Multi-Layer Perceptron (MLP) and a Convolutional Neural Network (CNN) based implementation.

Double DQN Architecture (DDQN)

The Double Deep Q-Network (DDQN) architecture is an extension of the standard DQN designed to mitigate the problem of overestimation of Q-values. In traditional DQN, the same Q-network is used for both selecting and evaluating actions, which can lead to optimistic value estimates and potentially cause the agent to learn suboptimal policies.

DDQN addresses this limitation by decoupling the action selection from the action evaluation. It introduces two separate Q-networks:

- **Online Q-Network:** This network is actively updated during training and is responsible for selecting the greedy action for the next state.
- **Target Q-Network:** This network is a periodically updated copy of the Online Q-Network (updated less frequently). It is used to evaluate the Q-value of the action chosen by the Online Q-Network.

The target Q-value calculation in DDQN is modified to prevent overestimation. Instead of using the same network for both selecting and evaluating the next action, the Online Q-Network selects the best action, while the Target Q-Network evaluates it. The modified target is given by:

$$q_{\text{ref}} = r_t + \lambda \cdot Q_{\text{target}}(s_{t+1}, \arg \max_{a'} Q_{\text{online}}(s_{t+1}, a'))$$

This formulation ensures that the target value is based on a more stable and less biased estimate, effectively reducing the positive bias commonly observed in standard DQN.

The overall training process of DDQN follows that of standard DQN, with the key difference being in the construction of the target Q-value. By avoiding the overestimation of action values, DDQN leads to more stable learning and can yield improved performance, especially in environments where value overestimation significantly affects policy quality.

CNN-DQN Architecture

The Convolutional Neural Network (CNN)-based Q-network represents a sophisticated approach built on the hypothesis that meaningful, localized patterns and spatial correlations exist within the network traffic feature set. Unlike the Multi-Layer Perceptron (MLP), which treats input as a flat feature vector, the CNN is specifically designed to exploit spatial structures, enabling automatic and hierarchical feature extraction.

Transforming Tabular Data into a Feature-Image

At the core of this method is the transformation of the one-dimensional feature vector into a two-dimensional *feature-image*. This enables the application of convolutional operations, framing the classification task as an image recognition problem.

Given a state vector with N features, it is reshaped into a grid of dimensions $H \times W$ such that $H \cdot W = N$. For instance, a 1D vector of 64 features can be reshaped into an 8×8 grid as follows:

- Features 1–8 form the first row.
- Features 9–16 form the second row.
- ...
- Features 57–64 form the eighth row.

Each cell in this 8×8 grid is treated as a pixel, where its value corresponds to a normalized feature value, resulting in a single-channel input image.

Feature Ordering Strategy

The order in which features are arranged in the grid critically impacts the CNN’s ability to detect meaningful patterns. Several strategies can be adopted:

- **Default Ordering:** Based on the output from the feature extraction pipeline.
- **Correlation-Based Ordering:** Positions highly correlated features adjacent to one another using a correlation matrix.
- **Domain Knowledge-Based Ordering:** Groups features by semantics, e.g., timing-based, size-based, or protocol-based.

These strategies aim to ensure that local neighborhoods within the feature-image encode relevant, structured information.

CNN Architecture

Input Normalization: A 2D Batch Normalization layer is applied to stabilize input variance.

First Convolutional Block (Low-Level Pattern Detection): • **Convolutional Layer:** 16 filters of size 3×3 .

- **Activation:** ReLU.
- **Max Pooling:** Reduces spatial resolution, enabling translation invariance.
- **Spatial Dropout:** Dropout rate of 0.25 applied to 2D feature maps.

Second Convolutional Block (High-Level Abstraction): • **Convolutional Layer:** 32 filters of size 3×3 .

- **Activation:** ReLU.
- **Dropout:** Spatial dropout applied again for regularization.

Fully Connected Head (Q-Value Estimation): • **Flattening:** Converts the feature map into a 1D vector.

- **Dense Layer:** Fully connected layer with 64 neurons, followed by BatchNorm and ReLU.
- **Dropout:** Dropout rate of 0.5 to prevent overfitting.
- **Output Layer:** Final linear layer mapping to the number of possible actions (e.g., 2 for binary classification), producing raw Q-values.

These Q-values are suitable for downstream optimization using a classification loss such as Cross-Entropy. Softmax may be applied implicitly, depending on the loss function implementation.

3.2.4 Conclusion

This chapter has presented a comprehensive exposition of the design of our DQN-based DDoS detection system, with a specific emphasis on adapting reinforcement learning principles to a supervised classification context.

We began by modeling the problem through formal definitions of environment states, action spaces, reward signals, and the agent’s learning objective. Particular attention was given to the construction of the reward function, which serves as a crucial link between reinforcement learning and supervised labels.

The chapter also detailed the preprocessing pipeline, encompassing dataset cleaning, feature selection, normalization, and class balancing. These steps were shown to be essential in ensuring the quality and consistency of the data fed into the learning models.

Next, we introduced the architectural foundations of our proposed system, starting with the Deep Q-Network (DQN), extending to the Double DQN (DDQN) architecture for reducing overestimation bias, and

culminating in the CNN-DQN design—leveraging convolutional layers to detect spatial patterns in structured network traffic data.

Together, these methodological components lay the groundwork for the implementation, training, and empirical evaluation phases discussed in the upcoming chapters, wherein the effectiveness of these approaches will be rigorously assessed in real-world DDoS detection scenarios.

Bibliography

- [1] IBM, *Security Fundamentals?*,
<https://www.linkedin.com/pulse/fundamentals-security-slamghana/>
- [2] CIA, *Principles of Information Security*.
Available at: <https://www.cleanpng.com/png-information-security-confidentiality-availability-1373941/>
- [3] Fidelis Security, *Understanding DoS and DDoS Attacks*,
<https://fidelissecurity.com/blog/dos-ddos-attacks>
- [4] Imperva, *Denial of Service (DoS) Attack*,
<https://www.imperva.com/learn/ddos/denial-of-service/>
- [5] NETSCOUT, *Slowloris Attack*,
<https://www.netscout.com/blog/asert/slowloris-attack>
- [6] MDPI, *Anomaly-based Detection of DDoS Attacks Using Machine Learning*,
<https://www.mdpi.com/2076-3417/10/3/1052>
- [7] IBM, *What is an Intrusion Detection System (IDS)?*,
<https://www.ibm.com/think/topics/intrusion-detection-system>
- [8] Wikipedia, *Intrusion Detection System*,
https://en.wikipedia.org/wiki/Intrusion_detection_system
- [9] Fortinet, *What is Intrusion Detection Systems (IDS)? How does it Work?*,
<https://www.fortinet.com/resources/cyberglossary/intrusion-detection-system>
- [10] TechTarget, *What Is an Intrusion Detection System (IDS)?*,
<https://www.techtarget.com/searchsecurity/definition/intrusion-detection-system>
- [11] TechTarget, *Dataset Definition*,
<https://www.unb.ca/cic/datasets/ddos-2019.html>