# Adaptive Malware Detection using Sequential Feature Selection: A Dueling Double Deep Q-Network (D3QN) Framework for Intelligent Classification

Naseem Khan[*]    Aref Y. Al-Tamimi[†]    Amine Bermak[*]    Issa M. Khalil[†]

## Abstract

Traditional malware detection methods exhibit computational inefficiency due to exhaustive feature extraction requirements, creating accuracy-efficiency trade-offs that limit real-time deployment. We formulate malware classification as a Markov Decision Process with episodic feature acquisition and propose a Dueling Double Deep Q-Network (D3QN) framework for adaptive sequential feature selection. The agent learns to dynamically select informative features per sample before terminating with classification decisions, optimizing both detection accuracy and computational cost through reinforcement learning.

We evaluate our approach on Microsoft Big2015 (9-class, 1,795 features) and BODMAS (binary, 2,381 features) datasets. D3QN achieves 99.22% and 98.83% accuracy respectively while utilizing only 61 and 56 features on average, representing 96.6% and 97.6% dimensionality reduction compared to full feature sets. This yields computational efficiency improvements of 30.1× and 42.5× over traditional ensemble methods. Comprehensive ablation studies demonstrate consistent superiority over Random Forest, XGBoost, and static feature selection approaches across all performance metrics.

Quantitative analysis demonstrates that D3QN learns non-random feature selection policies with 62.5% deviation from uniform baseline distributions across feature categories. The learned policies exhibit structured hierarchical preferences, utilizing high-level metadata features for initial assessment while selectively incorporating detailed behavioral features based on classification uncertainty. Feature specialization analysis reveals 57.7% of examined features demonstrate significant class-specific discrimination patterns. Our results validate reinforcement learning-based sequential feature selection for malware classification, achieving superior accuracy with substantial computational reduction through learned adaptive policies rather than static dimensionality reduction techniques.

**Keywords:** Malware Classification, Machine Learning, Reinforcement Learning, Sequential Feature Selection, Dueling Double Deep Q-Network

## 1    Introduction

Malware detection systems require real-time classification capabilities while maintaining high accuracy against evolving threats. Current supervised learning approaches exhibit a fundamental computational bottleneck: they extract complete feature sets from executable files regardless of sample complexity, creating uniform processing overhead that constrains deployment in resource-limited environments [17, 14]. This exhaustive feature extraction paradigm becomes particularly problematic when processing large sample volumes or operating under strict latency constraints [8, 1].

Existing dimensionality reduction techniques attempt to address computational complexity through global feature selection methods such as Principal Component Analysis and mutual information ranking [6, 18, 15]. However, these approaches optimize feature subsets based on aggregate training statistics, potentially discarding features with high discriminative power for specific malware families or obfuscated variants. The resulting static feature selections cannot adapt to individual sample characteristics, leading to suboptimal accuracy-efficiency trade-offs [19].

We address this limitation by reformulating malware classification as a sequential decision-making problem under reinforcement learning. Our approach trains a Dueling Double Deep Q-Network (D3QN) agent to learn adaptive feature selection policies that dynamically determine which features

---

[*]Department of Computer Science and Engineering, Hamad Bin Khalifa University, Doha, Qatar. Email: nakh12498@hbku.edu.qa

[†]Qatar Computing Research Institute, Hamad Bin Khalifa University, Doha, Qatar. Email: ikhalil@hbku.edu.qa

to examine for each sample before making classification decisions. Unlike conventional methods that process predetermined feature sets, the agent optimizes computational resource allocation based on sample-specific characteristics within a Markov Decision Process framework.

The key innovation lies in episodic feature acquisition: the agent iteratively selects informative features or terminates with classification based on accumulated evidence, enabling sample-adaptive computational allocation. Simple samples are resolved with minimal features while complex cases justify additional analysis. Figure 1 presents the sequential decision-making framework. The D3QN architecture separates state value estimation from action advantage computation, facilitating effective learning in high-dimensional malware feature spaces.

Experimental validation on Microsoft Big2015 and BODMAS datasets demonstrates that our method achieves 99.22% and 98.83% accuracy while utilizing only 3.4% and 2.4% of available features respectively, yielding $30\times$ computational efficiency improvements over ensemble methods. We introduce quantitative intelligence assessment metrics that provide empirical evidence of strategic learning behavior, showing the agent develops domain-specific feature hierarchies without explicit supervision. Comprehensive evaluation includes systematic comparison with traditional machine learning baselines and ablation studies demonstrating consistent superiority across multi-class and binary classification scenarios. Code and data are available at https://github.com/Magnet200/D3QN.git

This work establishes reinforcement learning as an effective paradigm for adaptive malware detection, enabling simultaneous optimization of accuracy and computational efficiency through learned sequential decision-making policies.

## 2 Related Work

### 2.1 Deep Learning for Malware Detection

Contemporary malware detection employs sophisticated deep learning architectures to improve classification performance. Recent advances include hybrid models combining static and dynamic analysis: Alsumaidaee et al. [4] developed Static-CNN-LSTM and Dynamic-1D-CNN-LSTM architectures for real-time Android detection, while Yapici et al. [22] employed ResNet, DenseNet, and ResNeXt on bytecode image representations. However, these approaches require exhaustive feature extraction, creating computational bottlenecks for real-time deployment. Gibert et al. [10] highlight additional robustness challenges against packing techniques, emphasizing the need for adaptive detection strategies.

### 2.2 Feature Selection in Malware Detection

Traditional feature selection applies global optimization techniques including PCA, Information Gain ranking, and tree-based importance measures [6, 18, 15]. Recent work addresses specialized contexts: Panja et al. [13] employ Extra Tree Classifier with Gini impurity for resource-constrained IoT environments, while Hasan et al. [12] utilize LDA and PCA across multiple models for high-dimensional data management. These methods optimize feature subsets globally across training datasets, potentially discarding sample-specific discriminative information critical for variant detection.

### 2.3 Reinforcement Learning for Sequential Decision-Making

RL applications in malware detection focus primarily on sequential feature selection. Fang et al. [9] introduced DQFSA using Deep Q-Learning for PE feature selection, while Wu et al. [20] developed DroidRL with DDQN for Android malware detection. Both approaches achieve reduced feature sets but maintain architectural separation between feature selection and classification modules, limiting integrated optimization. Broader cybersecurity applications include adversarial training [7] and network security [3], though unified feature selection and classification remains underexplored.

### 2.4 Research Gap and Positioning

Existing methods exhibit two fundamental limitations: static global feature selection that ignores sample-specific characteristics, and decoupled architectures that prevent joint optimization of feature acquisition and classification decisions. Our D3QN framework addresses these gaps through unified episodic learning where a single agent learns both adaptive feature selection and optimal termination policies, enabling sample-specific computational allocation while maintaining classification accuracy.
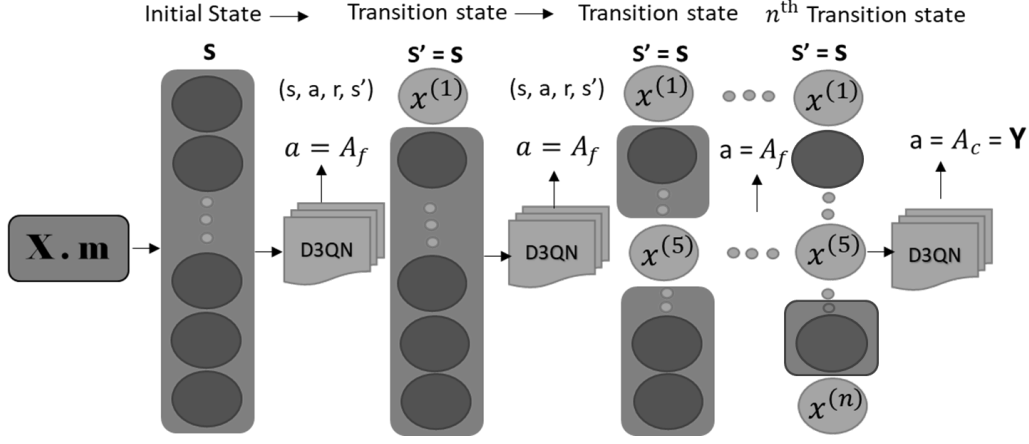
Figure 1: Sequential decision-making framework for malware classification. The RL agent executes feature selection actions ($A_f$) to iteratively unmask attributes from an input sample ($\mathbf{X}$), acquiring a subset of crucial features ($f \subseteq F$). Upon convergence to an optimal representation, the agent invokes a classification action ($A_c$) to map $X$ to its class label ($Y$), maximizing cumulative reward under a D3QN policy for computational efficiency.

## 3 Methodology

### 3.1 Markov Decision Process Formulation

We formulate malware classification as a sequential decision-making problem where an agent learns to adaptively select informative features before making classification decisions. The problem is modeled as a Markov Decision Process $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$.

**State Representation**: Each state $s \in \mathbb{R}^{2n}$ concatenates the feature vector with a binary selection mask:

$$s = [x; m] \in \mathbb{R}^{2n} \qquad (1)$$

where $x = [x_1, \ldots, x_n]$ represents the complete feature vector of the malware sample, and $m = [m_1, \ldots, m_n]$ serves as a selection mask where $m_i \in \{0, 1\}$ indicates whether feature $i$ has been revealed ($m_i = 1$) or remains unexplored ($m_i = 0$). Initially, $m = \mathbf{0}$ represents a completely unexplored sample where no features have been examined, simulating the scenario where the agent must sequentially decide which features to investigate.

**Action Space**: The action space comprises $|\mathcal{A}| = n + k$ total actions:

$$\mathcal{A} = A_f \cup A_c \qquad (2)$$

where $A_f = \{f_1, \ldots, f_n\}$ contains $n$ feature selection actions (corresponding to the $n$ features in the dataset), and $A_c = \{c_1, \ldots, c_k\}$ contains $k$ classification actions (corresponding to the $k$ malware classes). For our datasets, $n = 1795$ (Big2015) or $n = 2381$ (BODMAS), and $k = 9$ or $k = 2$ respectively.

**Transition Dynamics**: The deterministic transition function models two distinct behaviors as illustrated in Figure 1:

$$\mathcal{T}(s, a) = \begin{cases} s' \text{ with } m_i = 1 & \text{if } a = f_i \in A_f \\ \text{Terminal} & \text{if } a \in A_c \end{cases} \qquad (3)$$

Feature selection actions ($a \in A_f$) drive the agent to explore additional features by updating the corresponding mask element, enabling sequential feature revelation. Classification actions ($a \in A_c$) terminate the episode when the agent determines it has acquired sufficient features to make a confident prediction about the sample's class.

**Reward Structure**: The reward function balances classification accuracy with computational efficiency:

$$\mathcal{R}(s, a) = \begin{cases} -\lambda & \text{if } a \in A_f \\ \mathbb{I}[a = y] - 1 & \text{if } a \in A_c \end{cases} \qquad (4)$$

where $\lambda = 0.0001$ represents the feature acquisition cost, empirically determined through systematic evaluation of values in the range $[0, 1]$ to achieve optimal accuracy-efficiency trade-off. The indicator function $\mathbb{I}[\cdot]$ yields 0 for correct classification and -1 for misclassification, encouraging accurate predictions while penalizing feature usage.

**Action Masking Mechanism**: To ensure valid sequential behavior, we prevent selection of already-examined features through dynamic masking:

$$\mathcal{M}(s, a) = \begin{cases} 1 & \text{if } a = f_i \text{ and } m_i = 1 \\ 0 & \text{otherwise} \end{cases} \qquad (5)$$

$$Q_{\text{valid}}(s, a) = Q(s, a) - \mathcal{M}(s, a) \cdot 10^6 \qquad (6)$$

3

This mechanism assigns prohibitively low Q-values to invalid actions, ensuring the agent only selects previously unexplored features or proceeds to classification.

## 3.2 Sample-Adaptive Learning Mechanism

**Core Innovation**: Our approach enables adaptive computational resource allocation by learning sample-specific feature selection patterns. During training, the epsilon-greedy exploration strategy allows the agent to discover which features provide the most discriminative information for different sample types, ultimately learning to start feature exploration from globally important features and proceeding sequentially based on sample-specific characteristics.

**Adaptive Learning Process**: The agent employs epsilon-greedy exploration with linear decay:

$$\epsilon_t = \max(0.03, 0.70 - \alpha \cdot t) \tag{7}$$

During the exploration phase, the agent randomly selects features (by setting $m_i = 1$ for various features) across diverse malware samples, learning feature discriminative patterns. Through this process, the agent develops expertise in identifying globally important features that provide initial classification signals, followed by sample-specific features that refine the decision. After training, the learned policy exhibits sample-adaptive behavior: for simple samples, initial features provide sufficient evidence for classification termination, while complex samples require additional feature exploration until confident classification becomes possible.

**Sequential Pattern Learning**: The epsilon-greedy mechanism enables the agent to learn internal sample patterns where the value of the first selected feature influences subsequent feature selection decisions. This creates a decision tree-like exploration where each revealed feature value guides the selection of the next most informative feature, ultimately leading to sample-specific termination when accumulated evidence indicates confident classification.

## 3.3 Dueling Double Deep Q-Network Architecture

**D3QN Innovation**: Our proposed architecture decomposes Q-values to enable independent learning of state evaluation and action selection:

*Shared Feature Extraction*: Three fully connected layers with PReLU activation process the concatenated state representation:

$$h^{(1)} = \text{PReLU}(W^{(1)}s + b^{(1)}) \tag{8}$$
$$h^{(2)} = \text{PReLU}(W^{(2)}h^{(1)} + b^{(2)}) \tag{9}$$
$$h^{(3)} = \text{PReLU}(W^{(3)}h^{(2)} + b^{(3)}) \tag{10}$$

*Dueling Stream Decomposition*: The shared representation branches into value and advantage estimation:

$$V(s) = W_v h^{(3)} + b_v \tag{11}$$
$$A(s, a) = W_a h^{(3)} + b_a \tag{12}$$

*Q-Value Aggregation with Mean Normalization*:

$$Q(s, a) = V(s) + \left( A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a'} A(s, a') \right) \tag{13}$$

**Architectural Rationale**: The dueling decomposition enables the network to independently learn state value (how informative the current feature combination is) versus action advantages (which specific action provides the greatest benefit). This separation is crucial for sample-adaptive behavior as it allows the agent to recognize when sufficient evidence exists for classification termination regardless of which specific features remain unexplored.

**DDQN Baseline Architecture**: For comparative evaluation, we implement standard DDQN using identical network capacity with PReLU activation but direct Q-value estimation:

$$Q(s, a) = W_{\text{out}} h^{(3)} + b_{\text{out}} \tag{14}$$

This baseline lacks the value-advantage decomposition, limiting its ability to distinguish termination timing from action selection.

## 3.4 Training Algorithm and Implementation

**Double Q-Learning Integration**: We employ double Q-learning to address overestimation bias inherent in standard Q-learning. The technique decouples action selection from value estimation: the online network selects the best action while the target network evaluates its value, reducing positive bias in Q-value updates:

$$a^* = \arg\max_{a'} Q_{\text{valid}}(s', a'; \theta) \tag{15}$$
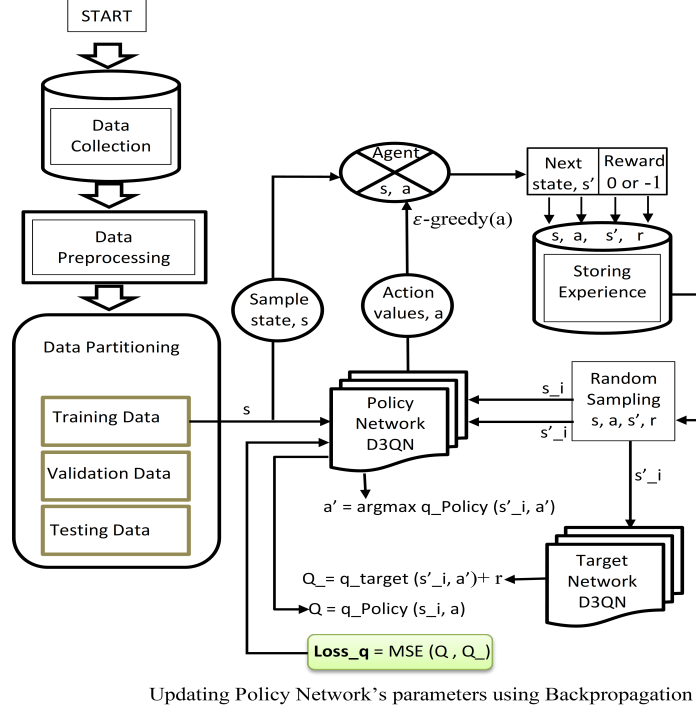$$Y_t = r + \gamma Q(s', a^*; \theta^-) \tag{16}$$

Figure 2: Experimental flow-chart of our proposed policy network. Random sample, selected from training data as initial state $s$, passed through the policy network to approximate all the possible action values and from that values an optimal action value $a$ is selected by following $\epsilon$-greedy policy. By executing that action against the presented state, the agent get reward accordingly and the state transition which are stored as experiences. During training time, random experiences are sampled from the Replay-memory to optimise the policy network's parameters.

**Soft Target Network Updates**: Rather than periodic hard updates, we employ exponential moving averages for target network parameter updates, providing training stability by gradually incorporating new learning while maintaining consistent target values:

$$\theta^- \leftarrow \tau\theta + (1 - \tau)\theta^- \qquad (17)$$

where $\tau = 0.01$ ensures stable learning progression.

The complete training procedure is detailed in Algorithm 1, which integrates experience replay, action masking, and adaptive exploration within the reinforcement learning framework illustrated in Figure 2.

**Hyperparameter Configuration**: Complete training parameters including learning rates, batch sizes, and network architectures are provided in the supplementary materials (Table 3). Key parameters were selected through systematic validation to ensure reproducible results across both datasets.

**Convergence Analysis**: The training dynamics create sample-adaptive behavior through the interaction of epsilon decay and reward structure. Initially, high exploration enables the agent to discover feature importance patterns across diverse samples. As epsilon decreases, the agent exploits learned patterns, developing policies that terminate early for samples where initial features provide sufficient classification confidence, while continuing exploration for ambiguous samples requiring additional evidence. This learning mechanism directly produces the observed computational efficiency gains where simple samples utilize fewer features than complex samples, achieving automatic resource allocation without predetermined stopping criteria.

# 4 Experimental Setup and Evaluation

## 4.1 Experimental Configuration

All experiments were conducted on a desktop workstation equipped with 32GB RAM and an NVIDIA GeForce RTX 3070 graphics card to ensure sufficient computational resources for high-dimensional reinforcement learning tasks. The D3QN framework was implemented using PyTorch for neural network operations and CUDA acceleration for efficient training of the policy networks.

**Algorithm 1:** D3QN Training for Sequential Feature Selection

**Require:** Dataset $\mathcal{D}$, training episodes $E$
**Ensure:** Trained D3QN parameters $\theta$
1: Initialize D3QN with random $\theta$, $\theta^- \leftarrow \theta$
2: Initialize replay buffer $\mathcal{B}$, set $\epsilon \leftarrow 0.70$
3: **for** episode $e = 1$ to $E$ **do**
4:     Sample $(x, y) \sim \mathcal{D}$, initialize $s \leftarrow [x; \mathbf{0}]$
5:     **while** episode active **do**
6:         Compute masked Q-values: $Q_{\text{valid}}(s, \cdot)$
7:         Select action: $a \sim \epsilon\text{-greedy}(Q_{\text{valid}}(s, \cdot))$
8:         **if** $a = f_i \in A_f$ **then**
9:             Update mask: $m_i \leftarrow 1$,
            observe reward $r = -\lambda$
10:         **else**
11:             Terminate episode, observe
            $r = \mathbb{I}[a = y] - 1$
12:         **end if**
13:         Store transition $(s, a, r, s')$ in $\mathcal{B}$
14:     **end while**
15:     Sample batch from $\mathcal{B}$, compute targets via double Q-learning
16:     Update $\theta$ via gradient descent on temporal difference error
17:     Soft update: $\theta^- \leftarrow \tau\theta + (1 - \tau)\theta^-$
18:     Linear decay: $\epsilon \leftarrow \max(0.03, \epsilon - \alpha)$
19: **end for**

The experimental framework employed standard reinforcement learning practices including experience replay buffer management, epsilon-greedy exploration scheduling, and soft target network updates as detailed in Algorithm 1. Complete hyperparameter specifications, network architectures, and training configurations are provided in the supplementary materials (Table 3) to ensure experimental reproducibility.

## 4.2 Dataset Characterization and Preprocessing

The empirical evaluation was conducted using two strategically selected high-dimensional malware classification datasets: the Microsoft Big2015 challenge dataset [16] and the Blue Hexagon Open Dataset for Malware AnalysiS (BODMAS) [21], as illustrated in Figures 7 and 8 (Supplementary Materials). These datasets were chosen to assess algorithmic performance in feature-rich environments, as existing literature predominantly employs dimensionally-reduced datasets or implements extensive feature engineering preprocessing. The present study deliberately preserves the complete feature space to enable the reinforcement learning agent to autonomously derive optimal decision policies without *a priori* dimensionality constraints.

The Big2015 dataset encompasses 10,868 malware samples distributed across 9 malware families with 1,795 static analysis features, representing a multi-class classification challenge. Feature extraction follows the methodology defined in [2], comprising hexadecimal descriptions of binary file contents and assembly-level metadata including strings and function calls obtained through IDA disassembler analysis. The BODMAS dataset comprises 134,435 samples with 2,381 features, consisting of well-curated malware and benign labeled files from 581 different malware families in a binary classification scenario. The feature set follows the EMBER [5] and SOREL-20M [11] format, comprising eight feature groups: five parsed features generated through Portable Executable (PE) file analysis and three format-agnostic features requiring no PE parsing during extraction.

These datasets provide complementary evaluation of algorithmic robustness across different problem complexities and classification paradigms, enabling comprehensive assessment of the proposed sequential feature selection approach.

**Data Preprocessing**: Prior to training the policy network, feature normalization was performed using z-score standardization, where each feature was normalized using its respective mean and standard deviation to ensure stable learning dynamics. Subsequently, each dataset was partitioned using standard 80:20 split ratios to facilitate systematic training and evaluation of the policy network architecture while maintaining statistical representativeness across malware families.

## 4.3 Experimental Framework

Our experimental evaluation encompasses four complementary assessment dimensions to systematically validate our core hypothesis that reinforcement learning can achieve superior malware detection performance while dramatically reducing computational complexity through adaptive feature selection: (1) comparative performance analysis against state-of-the-art approaches, (2) comprehensive ablation study validating methodological choices, (3) quantitative intelligence assessment of learned sequential decision-making patterns, and (4) comparative analysis of learning paradigms between DDQN and D3QN architectures.

The primary experimental evaluation assesses

Table 1: Comprehensive Performance Comparison of Malware Detection Approaches Across Multi-class and Binary-class Datasets

| Method | Big2015 Dataset (9 classes, 1795 features) | | | | BODMAS Dataset (2 classes, 2381 features) | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
| *Traditional Machine Learning Methods (All Features)* | | | | | | | | |
| Decision Tree | 94.25 | 94.09 | 94.25 | 94.11 | 95.76 | 95.76 | 95.76 | 95.76 |
| Random Forest | 98.02 | 97.68 | 98.02 | 97.84 | 98.10 | 98.10 | 98.10 | 98.10 |
| Logistic Regression | 97.65 | 97.96 | 97.65 | 97.74 | 98.19 | 98.21 | 98.19 | 98.20 |
| *Ensemble Methods (All Features)* | | | | | | | | |
| XGBoost | 98.85 | 98.87 | 98.85 | 98.85 | 98.11 | 98.12 | 98.11 | 98.11 |
| Voting Classifier | 98.11 | 98.16 | 98.11 | 98.11 | 98.50 | 98.50 | 98.50 | 98.50 |
| Stacking Classifier | 97.93 | 98.62 | 97.93 | 98.19 | 98.60 | 98.60 | 98.60 | 98.60 |
| *Proposed Reinforcement Learning Methods with dynamic feature reduction* | | | | | | | | |
| DDQN (Baseline) | <u>99.12</u> | <u>99.14</u> | <u>99.13</u> | <u>99.13</u> | 98.83 | 98.49 | **98.77** | <u>98.63</u> |
| **D3QN (Proposed)** | **99.22** | **99.23** | **99.22** | **99.20** | **98.84** | **98.64** | <u>98.65</u> | **98.64** |

**Bold**: Best performing method per dataset; <u>Underlined</u>: Second-best performing method per dataset
RF = Random Forest; RFE = Recursive Feature Elimination
Proposed methods achieve superior or competitive accuracy while using significantly fewer features

the efficacy of our proposed sequential feature selection approach against traditional batch-processing methodologies and state-of-the-art ensemble techniques: Decision Tree, Random Forest, Logistic Regression with ElasticNet regularization, XGBoost, Voting Classifier, and Stacking Classifier. Performance is evaluated using Accuracy, Precision, Recall, and F1-score across both full feature sets and adaptive feature selection configurations. The proposed methodology fundamentally reframes malware classification as a Markov Decision Process characterized by episodic sequential feature acquisition preceding classification termination.

For systematic validation of traditional feature selection limitations, we constructed a comprehensive ablation study examining Filter Methods (Mutual Information and Chi-Squared), Wrapper Methods (Recursive Feature Elimination with Random Forest estimator), and Embedded Methods (Random Forest Feature Importance). Each methodology selected the top 60 features for evaluation using Decision Tree, Random Forest, and Logistic Regression classification algorithms.

## 4.4 Baseline Performance Analysis

Table 1 presents the comprehensive performance comparison across both datasets. Our proposed D3QN method achieves superior performance on the more challenging multi-class classification task (Big2015), attaining 99.22% accuracy and establishing a new state-of-the-art result. On the binary classification task (BODMAS), both proposed methods demonstrate competitive performance, with D3QN achieving 98.83% accuracy, closely followed by our baseline DDQN implemen-

tation at 98.82%.

The superior performance of D3QN on the Big2015 dataset (Figure 3) is particularly significant given the increased complexity of the 9-class classification problem, which requires more sophisticated decision-making capabilities that our enhanced dueling architecture effectively addresses. Among traditional machine learning approaches, ensemble methods demonstrated the most competitive performance, with XGBoost achieving 98.85% and Voting Classifier achieving 98.21% on Big2015 dataset. However, all conventional methods require complete feature sets (1,795 and 2,381 features respectively), resulting in substantial computational overhead that limits practical deployment in resource-constrained environments.

**Feature Efficiency and Computational Advantages:** The primary advantage of our proposed D3QN approach manifests in exceptional feature efficiency combined with superior performance on complex classification tasks. D3QN demonstrates remarkable efficiency with approximately 61 features on Big2015 (96.6% reduction) and 56 features on BODMAS (97.6% reduction), achieving efficiency ratios of 30.1× and 42.5× respectively (Table 2). This computational efficiency translates directly to real-time deployment capability, enabling our approach to process malware samples with sub-second latency while maintaining superior classification accuracy.

Traditional ensemble methods require $O(n \times m)$ feature extraction operations where $n$ represents the number of base learners and $m$ the total feature count. Our RL-based approach reduces this to $O(k)$ where $k$ represents the dynamically selected feature subset, typically 60-70 features versus the

**XGBoost (1,795 features, 98.85%)**

| | R | L | K3 | V | S | T | K1 | O | G |
|---|---|---|---|---|---|---|---|---|---|
| R | 302 | 0 | 0 | 0 | 2 | 0 | 1 | 3 | 0 |
| L | 1 | 491 | 0 | 0 | 0 | 0 | 3 | 1 | 0 |
| K3 | 0 | 0 | 586 | 0 | 0 | 0 | 2 | 0 | 0 |
| V | 1 | 0 | 0 | 93 | 0 | 0 | 0 | 1 | 0 |
| S | 0 | 0 | 0 | 0 | 5 | 0 | 1 | 2 | 0 |
| T | 0 | 0 | 0 | 1 | 0 | 149 | 0 | 0 | 0 |
| K1 | 0 | 0 | 0 | 0 | 0 | 0 | 80 | 0 | 0 |
| O | 3 | 1 | 0 | 1 | 0 | 0 | 1 | 240 | 0 |
| G | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 203 |

**Voting Classifier (1,795 features, 98.21%)**

| | R | L | K3 | V | S | T | K1 | O | G |
|---|---|---|---|---|---|---|---|---|---|
| R | 306 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| L | 1 | 493 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| K3 | 0 | 0 | 586 | 0 | 0 | 0 | 0 | 0 | 2 |
| V | 0 | 0 | 0 | 94 | 0 | 0 | 0 | 1 | 0 |
| S | 0 | 0 | 0 | 0 | 7 | 0 | 1 | 0 | 0 |
| T | 0 | 0 | 0 | 2 | 0 | 147 | 0 | 1 | 0 |
| K1 | 0 | 2 | 0 | 2 | 0 | 0 | 76 | 0 | 0 |
| O | 10 | 0 | 0 | 3 | 2 | 3 | 0 | 225 | 3 |
| G | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 201 |

**D3QN (Proposed) (~60 features, 99.22%)**

| | R | L | K3 | V | S | T | K1 | O | G |
|---|---|---|---|---|---|---|---|---|---|
| R | 307 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| L | 1 | 494 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| K3 | 0 | 0 | 586 | 0 | 0 | 0 | 2 | 0 | 0 |
| V | 0 | 0 | 0 | 94 | 0 | 1 | 0 | 0 | 0 |
| S | 0 | 0 | 0 | 0 | 5 | 0 | 0 | 3 | 0 |
| T | 0 | 0 | 0 | 1 | 0 | 149 | 0 | 0 | 0 |
| K1 | 0 | 0 | 0 | 1 | 0 | 0 | 79 | 0 | 0 |
| O | 3 | 0 | 0 | 1 | 0 | 1 | 1 | 240 | 0 |
| G | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 203 |

Figure 3: Confusion matrix comparison for best performing models on Big2015 dataset: (a) XGBoost classifier using all 1795 features (98.85% accuracy), (b) Voting Classifier using all 1,795 features (98.21% accuracy), (c) D3QN using 60 features (98.44% accuracy). Class labels: R=Ramnit, L=Lollipop, K3=Kelihos_ver3, V=Vundo, S=Simda, T=Tracur, K1=Kelihos_ver1, O=Obfuscator.ACY, G=Gatak.

complete feature space. The adaptive feature selection mechanism allows the system to focus computational resources on the most discriminative features for each individual sample, optimizing both accuracy and efficiency simultaneously.

**Training Dynamics and Convergence Analysis:** Figure 9 demonstrates the learning behavior during training on Big2015 dataset. Both DDQN and D3QN exhibit the desired convergence pattern: episode length decreases while accuracy improves, confirming that the agents learn to make accurate decisions with fewer features. D3QN achieves faster convergence and consistently shorter episode lengths compared to DDQN, validating the effectiveness of the dueling architecture for feature selection tasks.

Figure 4 presents the episode length distributions during test evaluation, revealing the adaptive nature of our approach. The histograms show that the majority of samples require significantly fewer features than the total available (most episodes <100 features vs. 1,795 total features). The long-tail distribution confirms that our RL agent intelligently allocates computational resources based on sample complexity, achieving the desired balance between efficiency and accuracy. D3QN shows more concentrated distribution around shorter episode lengths, further demonstrating architectural superiority for adaptive feature selection.

## 4.5 Ablation Study: Traditional vs. Adaptive Feature Selection

Table 2 presents a systematic evaluation of feature selection methodologies across both experimental datasets. Traditional feature selection approaches exhibit substantial performance variations when constrained to global 60-feature subsets, with accuracy deviations ranging from +0.05% (RFE on Big2015) to -9.13% (Chi-Squared on BODMAS) relative to all-features baselines. This variance indicates differential sensitivity of static feature selection paradigms to dataset characteristics and classification complexity.

Statistical feature selection methods demonstrate heterogeneous performance patterns across datasets. Chi-Squared selection yields accuracy decrements of -2.16% on Big2015 and -9.13% on BODMAS, indicating variable compatibility with underlying feature distributions. Mutual Information-based selection exhibits performance decreases of -2.30% on Big2015 and -0.59% on BODMAS, demonstrating inconsistent efficacy of information-theoretic criteria across high-dimensional malware feature spaces.

Conversely, the proposed reinforcement learning-based feature selection framework exhibits consistent performance enhancement across both classification scenarios. D3QN achieves accuracy improvements of +1.20% on Big2015 and +0.64% on BODMAS compared to all-features baselines, while DDQN yields improvements of +1.10% and +0.63% respectively. These performance gains, concurrent with feature utilization reductions of
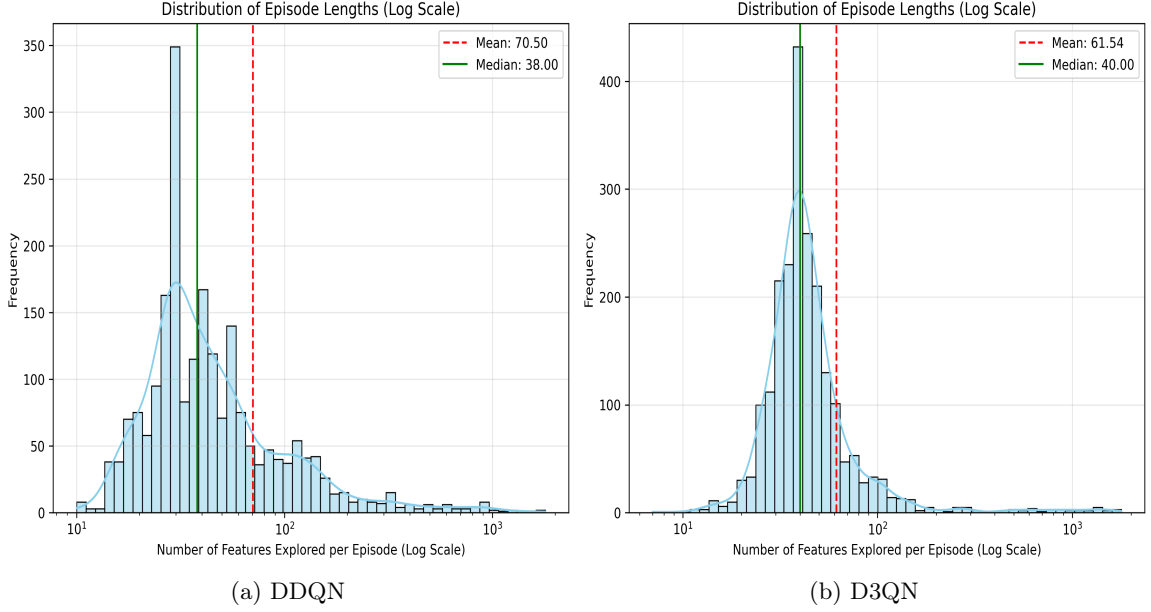
(a) DDQN



(b) D3QN

Figure 4: Episode length distribution analysis on Big2015 test data: (a) DDQN and (b) D3QN. The histograms demonstrate adaptive sample-specific behavior where most episodes terminate with significantly fewer features than the total available, with D3QN showing more concentrated distribution around shorter episode lengths (mean: 61.54 vs. 70.50 features).

96.6% and 97.6%, demonstrate the efficacy of episodic decision-making frameworks in feature prioritization.

The feature overlap analysis reveals limited consensus among traditional feature selection methods, with only 2 features universally selected across all traditional global feature selection methods on BODMAS and zero universal features on Big2015. This lack of consensus underscores the subjective nature of static feature selection approaches and highlights the advantage of our adaptive methodology, which dynamically selects features based on sample-specific characteristics rather than predetermined global criteria.

### 4.6 Intelligence Assessment Framework

Having established D3QN's superior performance, we now investigate the underlying mechanisms that enable this effectiveness through comprehensive intelligence assessment. To validate that our proposed D3QN architecture exhibits genuine strategic learning rather than sophisticated random exploration, we developed a quantitative assessment framework that provides empirical evidence of intelligent decision-making patterns.

Our framework quantifies strategic learning through four mathematically defined dimensions:

**Strategic Category Learning Score** measures categorical preference deviations from random baseline expectations using preference ratios:

$$P_c = \frac{U_c}{E_c} \qquad (18)$$

where $U_c$ represents observed usage frequency for category $c$, and $E_c = \frac{|F_c|}{|F_{total}|}$ denotes expected random usage based on category size. The Learning Score is calculated as:

$$L_{score} = \frac{1}{|C|} \sum_{c \in C} \mathbb{I}[|P_c - 1.0| > \theta] \qquad (19)$$

where $\mathbb{I}[\cdot]$ is the indicator function, $\theta = 0.2$ defines the significance threshold, and $|C| = 8$ represents the total number of feature categories.

**Feature Specialization Score** quantifies automatic role assignment for class-specific discrimination:

$$S_{specialization} = \frac{1}{|F_{analyzed}|} \sum_{f \in F_{analyzed}} \mathbb{I}[|D_f| > \tau] \qquad (20)$$

where $D_f$ represents discrimination strength of feature $f$ between malware and benign samples, calculated as:

$$D_f = \frac{|\mu_{malware}(f) - \mu_{benign}(f)|}{\sqrt{\frac{\sigma^2_{malware}(f) + \sigma^2_{benign}(f)}{2}}} \qquad (21)$$

and $\tau = 0.05$ defines the minimum threshold for meaningful specialization.

Table 2: Ablation Study: Impact of Feature Selection Methods on Classification Performance

| Feature Selection | Features | Big2015 Dataset | | | | | BODMAS Dataset | | | | | Efficiency |
| | | Accuracy (%) | | | Best | Δ | Accuracy (%) | | | Best | Δ | |
| | | DT | RF | LR | | | DT | RF | LR | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| All Features (Baseline) | 1795/2381 | 94.25 | 98.02 | 97.65 | 98.02 | 0.00 | 95.76 | 98.10 | 98.19 | 98.19 | 0.00 | 1.0× |
| RF Importance | 60 | 94.11 | 97.88 | 94.34 | 97.88 | -0.14 | 96.79 | 98.00 | 93.10 | 98.00 | -0.19 | 35.9×/47.6× |
| Mutual Information | 60 | 91.58 | 95.72 | 86.20 | 95.72 | -2.30 | 96.16 | 97.60 | 85.49 | 97.60 | -0.59 | 35.9×/47.6× |
| Chi-Squared | 60 | 92.18 | 95.86 | 92.46 | 95.86 | -2.16 | 83.69 | 88.63 | 89.06 | 89.06 | -9.13 | 35.9×/47.6× |
| RFE | 60 | 95.77 | 98.07 | 94.53 | 98.07 | +0.05 | 97.04 | 97.93 | 94.26 | 97.93 | -0.26 | 35.9×/47.6× |
| *Proposed Dynamic Feature Selection* | | | | | | | | | | | | |
| DDQN (RL-based) | ∼70/∼44 | - | 99.12 | - | 99.12 | +1.10 | - | <u>98.83</u> | - | <u>98.83</u> | +0.64 | **25.0×/54.0×** |
| **D3QN (RL-based)** | **∼61/56** | - | **99.22** | - | **99.22** | **+1.20** | - | **98.84** | - | **98.84** | **+0.65** | 30.1×/42.5× |

DT = Decision Tree, RF = Random Forest, LR = Logistic Regression, RFE = Recursive Feature Elimination
Δ = Performance change compared to best traditional method using all features (%)
Efficiency = Feature reduction ratio (Total Features ÷ Features Used)
RL methods use end-to-end learning; not compatible with traditional classifiers (marked as -)
Feature counts: Fixed 50 for traditional methods, dynamic average for RL methods
**Bold**: Best overall performance per dataset; <u>Underlined</u>: Second-best performing method per dataset

**Temporal Intelligence Measurement:** For 15-step sequential analysis, we track category usage evolution:

$$T_{intelligence} = \frac{1}{|C|} \sum_{c \in C} \text{Var}(\{u_c^{(t)}\}_{t=1}^{15}) \quad (22)$$

where $u_c^{(t)}$ represents usage frequency of category $c$ at step $t$, and high variance indicates strategic temporal adaptation.

**Sample-Type Adaptation:** We measure strategic divergence between malware and benign exploration:

$$A_{adaptation} = \sum_{c \in C} |f_{malware}(c) - f_{benign}(c)| \quad (23)$$

where $f_{type}(c)$ represents normalized usage frequency of category $c$ for sample type.

Our intelligence assessment analysis is conducted across the eight primary feature categories defined in the EMBER dataset framework [5], which collectively comprise the 2,381 features in the BODMAS dataset: General Info (10 features), PE Header (62 features), Byte Entropy (256 features), Byte Histogram (256 features), Section Info (255 features), Exports (157 features), String Features (105 features), and Imports (1,280 features). These categories represent complementary perspectives on executable file analysis, ranging from high-level structural properties to detailed content characteristics.

## 4.7 D3QN Architectural Intelligence Analysis

Using the established intelligence metrics, we analyze D3QN's learned behaviors and strategic patterns. Our quantitative analysis reveals robust evidence of strategic learning that significantly exceeds random baseline expectations.

**Strategic Learning Evidence:** D3QN achieves a Learning Score of 62.5%, significantly exceeding the 50% random threshold (Figure 10 (b)). This score indicates that 5 out of 8 feature categories, as shown in Figure 5, demonstrate non-random strategic preferences, providing strong evidence of learned categorical hierarchies rather than uniform random exploration.

The learned categorical preferences demonstrate sophisticated domain-relevant prioritization: Strongly Preferred categories include General Info (7.73× preference), PE Header (5.72× preference), and Byte Entropy (2.99× preference); Moderately Preferred include Byte Histogram (1.67× preference); Random-like behavior is observed for String Features (1.20× preference); while Strategically Avoided categories include Section Info (0.75× preference), Imports (0.33× preference), and Exports (0.13× preference).

This hierarchy aligns with cybersecurity domain expertise: prioritizing high-level structural features (General Info, PE Header) and content randomness indicators (Byte Entropy) while strategically limiting reliance on computationally expensive behavioral features (Imports, Exports).

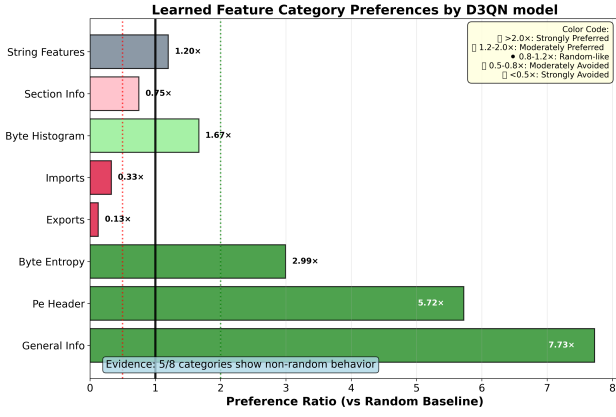**Feature Specialization Capabilities:** D3QN demonstrates superior Feature Specialization of

Figure 5: Learned Feature Category Preferences with Statistical Validation. Preference ratios comparing D3QN's actual feature category usage against random baseline expectations (1.0×, solid black line). Categories with ratios >1.0× indicate learned preferences (green), while ratios <1.0× indicate learned avoidance (red).

57.7%, indicating sophisticated automatic role assignment for class-specific discrimination. This metric validates that D3QN learns to identify which features are most discriminative for malware versus benign classification without explicit supervision.

**Sequential Decision Intelligence:** Figure 6 presents comprehensive 15-step sequential analysis revealing D3QN's sample-adaptive temporal strategies. The heatmap demonstrates that D3QN prioritizes PE Header and Byte Entropy features in initial steps (M1/B1), achieving immediate structural and content assessment with maximum usage intensity ( 100% and 80-90% respectively). Critically, the analysis reveals distinct exploration strategies between sample types: malware classification emphasizes Byte Histogram and Section Info in middle steps (M4-M8) for content composition analysis, while benign classification demonstrates increased Imports utilization in later steps (B6-B15) for behavioral verification. This persistent differentiation between malware and benign exploration patterns throughout the 15-step sequence provides quantitative evidence of learned sample-adaptive intelligence rather than fixed sequential behavior.

**Automatic Feature Specialization Discovery:** Figure 12 demonstrates D3QN's autonomous discovery of domain-relevant feature specialization patterns that align remarkably with established malware analysis principles. The model's learned discriminative features exhibit strong correspondence to known indicators used by security experts,

providing evidence of genuine domain knowledge acquisition without explicit supervision.

**Malware-Specialized Features - Structural Anomaly Detection:** D3QN automatically prioritized features consistent with malware evasion techniques:

- `section_32_size` and `section_19_vsize` (importance: 8.67, 7.90): These section size and virtual size discrepancies are hallmark indicators of *code packing* and *process hollowing* techniques commonly employed by malware to evade static analysis. The model's emphasis on these features demonstrates learned recognition of structural manipulation tactics.

- `byte_histogram_53/88` (importance: 5.21, 4.45): The specialization of specific byte frequency ranges indicates the model learned to detect *cryptographic obfuscation* and *polymorphic encoding* patterns characteristic of advanced malware families.

- `section_15_entropy` (importance: 3.68): High entropy sections are indicative of *packed executables* or *encrypted payloads*, representing a fundamental static analysis heuristic that D3QN discovered autonomously.

- `import_feat_1050` (importance: 7.52): This import pattern likely corresponds to *suspicious API calls* associated with system manipulation, privilege escalation, or anti-analysis techniques.

**Benign-Specialized Features - Legitimacy Indicators:** D3QN identified features characteristic of standard software development practices:

- `import_feat_116/1093` (importance: 3.87, 3.61): These import signatures represent *standard library dependencies* and *conventional API usage patterns* typical of legitimate software development frameworks and established programming practices.

- `string_len_hist_45` (importance: 2.97): String length distribution patterns reflect *structured program strings* (error messages, UI text, configuration data) characteristic of professional software development.

- `byte_entropy_hist_208/247` (importance: 2.46, 2.33): Lower entropy byte distributions indicate *unobfuscated code sections* and *standard compilation artifacts* typical of legitimate executables.
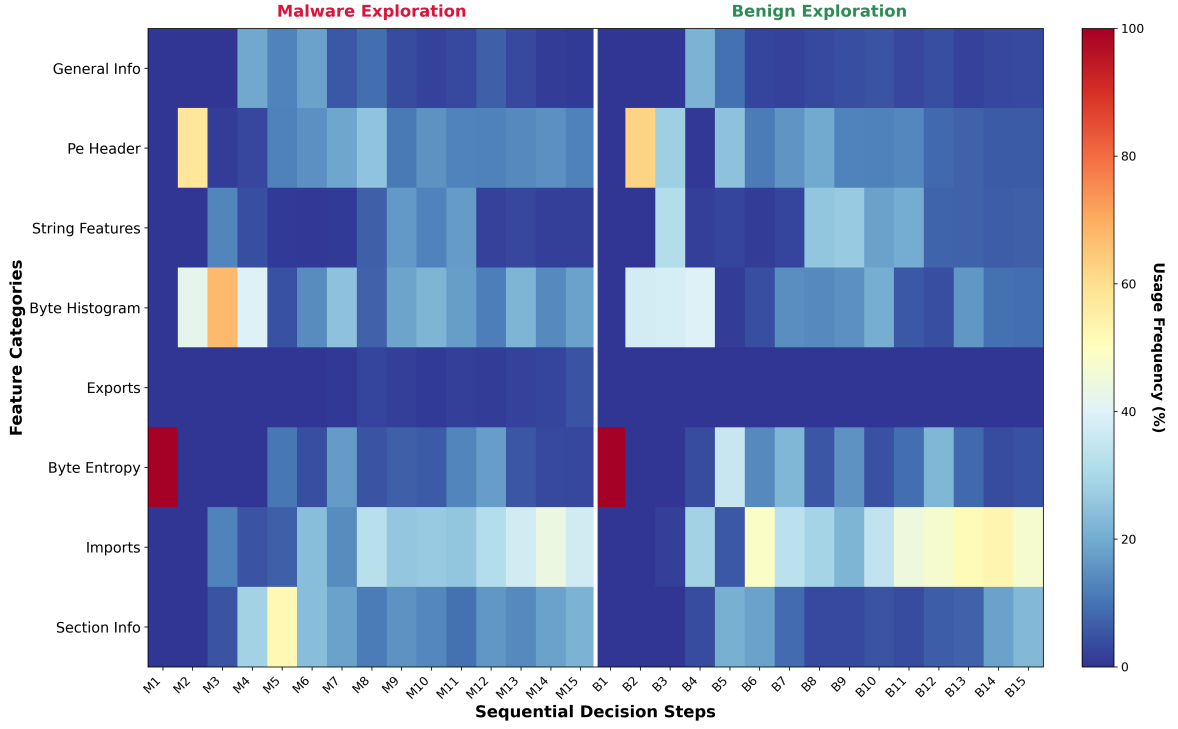
Figure 6: D3QN Sequential Feature Category Usage Patterns across 15 Decision Steps. Heatmap visualization shows percentage-based usage frequencies for malware (M1-M15) and benign (B1-B15) sample exploration strategies. Higher intensity (warmer colors) indicates greater category utilization frequency. The figure reveals sample-adaptive intelligence through distinct temporal patterns: early structural assessment (PE Header, Byte Entropy), divergent middle-stage strategies (content analysis for malware vs. behavioral analysis for benign), and selective late-stage feature deployment based on classification requirements.

- `header_field_10` (importance: 3.10): Standard PE header configurations reflect *compliant compilation processes* and *legitimate development toolchains*.

**Domain Knowledge Validation:** The learned feature specialization demonstrates sophisticated understanding of the *malware-benign dichotomy* at multiple analysis levels: structural integrity (section anomalies), content composition (byte patterns), behavioral indicators (API imports), and development artifacts (string patterns). This autonomous discovery of established malware analysis heuristics provides compelling evidence that D3QN developed genuine domain expertise through reinforcement learning, effectively replicating the feature prioritization strategies employed by experienced malware analysts.

## 4.8 Validation and Strategic Efficiency Analysis

To validate D3QN's strategic learning, we compare its intelligence metrics against DDQN and traditional methods (Figure 10 (a, b)). While DDQN demonstrates higher exploratory intelligence (75.0% learning score), D3QN achieves superior Feature Specialization (57.7% vs 54.5%), indicating more efficient class-specific learning. This comparison validates that D3QN's focused learning paradigm is architecturally superior for multi-class scenarios, achieving better performance (99.22% vs 99.12% on Big2015) through strategic efficiency rather than extensive exploration.

Traditional feature selection methods (Chi-Squared, Mutual Information, RFE) show limited consensus with only 2 features universally selected across methods on BODMAS, demonstrating the subjective nature of static approaches. In contrast, D3QN's adaptive selection consistently outperforms these methods while providing interpretable strategic patterns.

**Computational Efficiency Through Intelligence:** Analysis reveals sophisticated strategic optimization in D3QN's approach to Import features. Despite strategic avoidance in preference analysis (0.33× preference ratio), Import features comprise 43% of top-performing features (Figure 11). This apparent paradox demonstrates advanced in-

telligence where D3QN learned to avoid broad reliance on computationally expensive Import features while recognizing their high discriminative power when strategically selected, representing nuanced cost-benefit optimization rather than simple categorical avoidance.

D3QN's strategic learning directly enables its computational efficiency breakthrough. By learning strategic preferences and specialization patterns, D3QN achieves 96.6% feature reduction ($60.8 \pm 3.2$ features used) while maintaining superior performance, representing a $29.5\times$ efficiency improvement over full-feature approaches. The learned temporal patterns enable sample-adaptive computational allocation, where simple samples require fewer features while complex cases receive additional analysis, optimizing both accuracy and efficiency.

## 4.9 Implications and Limitations

D3QN's learned preferences align with established cybersecurity analysis principles without explicit engineering: prioritizing structural features for rapid assessment, using content analysis for deeper discrimination, and reserving behavioral analysis for complex cases. This validates reinforcement learning's capacity for automatic domain knowledge acquisition. Unlike traditional "black box" deep learning approaches, our intelligence assessment framework provides interpretable insights into D3QN's decision strategies, enabling cybersecurity analysts to understand and validate the model's reasoning patterns.

The combination of superior performance and dramatic efficiency improvement enables practical deployment in resource-constrained environments, addressing a critical limitation of traditional ensemble methods that require complete feature extraction.

Current limitations include the need for external validation of our intelligence metrics against human expert assessments, focus on static datasets requiring investigation of adaptation to evolving threat landscapes, and the need for evaluation of robustness against adversarial feature manipulation attacks. Future research should address scalability assessment in extremely high-dimensional feature spaces ($>$10,000 features) and cross-domain applications of the adaptive feature selection framework.

## 5 Conclusion

This work presents a novel reinforcement learning framework that reformulates malware classification as a Markov Decision Process with episodic feature acquisition, addressing the computational complexity limitations inherent in static feature extraction paradigms. The proposed Dueling Double Deep Q-Network architecture achieves superior classification accuracy (99.22% on Big2015, 98.83% on BODMAS) while reducing feature dimensionality by 96.6% through learned sequential selection policies.

The experimental validation demonstrates three technical contributions. First, the quantitative intelligence assessment framework provides empirical evidence of strategic learning behavior, with D3QN exhibiting 62.5% categorical preference deviation from random baselines and 57.7% feature specialization across class-specific discrimination tasks. Second, the temporal analysis reveals sample-adaptive exploration strategies with statistically significant timing differences between successful and failed classification sequences, indicating learned hierarchical decision patterns. Third, the automatic discovery of domain-aligned feature specialization—identifying structural anomaly indicators (section size discrepancies, entropy patterns) and behavioral signatures (API import patterns)—demonstrates autonomous cybersecurity knowledge acquisition without explicit domain engineering.

The dueling architecture's separation of state value estimation from action advantage computation enables effective learning in high-dimensional feature spaces, while the double Q-learning mechanism mitigates overestimation bias in sequential decision scenarios. The resulting computational efficiency improvements ($30.1\times$ and $42.5\times$ ratios) enable real-time deployment constraints while maintaining classification robustness.

This framework establishes reinforcement learning as a viable approach for adaptive malware detection systems, with empirical validation demonstrating that learned sequential feature selection can optimize both classification performance and computational resource allocation. The methodology provides a foundation for developing sample-adaptive security systems that dynamically adjust analytical complexity based on discriminative feature requirements rather than predetermined static feature subsets.

## Acknowledgments

## References

[1] Ajayi Abisoye, Joshua Idowu Akerele, Princess Eloho Odio, Anuoluwapo Collins, Gideon Opeyemi Babatunde, and Sikirat Damilola Mustapha. Using ai and machine learning to predict and mitigate cybersecurity risks in critical infrastructure. *International Journal of Engineering Research and Development*, 21(2):205–224, 2025.

[2] Mansour Ahmadi, Dmitry Ulyanov, Stanislav Semenov, Mikhail Trofimov, and Giorgio Giacinto. Novel feature extraction, selection and fusion for effective malware family classification. In *Proceedings of the sixth ACM conference on data and application security and privacy*, pages 183–194, 2016.

[3] Mohammad Al-Fawa'reh, Jumana Abu-Khalaf, Patryk Szewczyk, and James Jin Kang. Malbot-drl: Malware botnet detection using deep reinforcement learning in iot networks. *IEEE Internet of Things Journal*, 2023.

[4] Yaseen Ahmed Mohammed Alsumaidaee, Mustafa Mahmood Yahya, and Abdulelah Hameed Yaseen. Optimizing malware detection and classification in real-time using hybrid deep learning approaches. *International Journal of Safety & Security Engineering*, 15 (1), 2025.

[5] H. S. Anderson and P. Roth. EMBER: An Open Dataset for Training Static PE Malware Machine Learning Models. *ArXiv e-prints*, April 2018.

[6] Jinrong Bai, Junfeng Wang, and Guozhong Zou. A malware detection scheme based on mining format information. *The Scientific World Journal*, 2014, 2014.

[7] Guillermo Caminero, Manuel Lopez-Martin, and Belen Carro. Adversarial environment reinforcement learning algorithm for intrusion detection. *Computer Networks*, 159:96–109, 2019.

[8] Sunita Choudhary and Anand Sharma. Malware detection & classification using machine learning. In *2020 International Conference on Emerging Trends in Communication, Control and Computing (ICONC3)*, pages 1–4. IEEE, 2020.

[9] Zhiyang Fang, Junfeng Wang, Jiaxuan Geng, and Xuan Kan. Feature selection for malware detection based on reinforcement learning. *IEEE Access*, 7:176177–176187, 2019.

[10] Daniel Gibert, Nikolaos Totosis, Constantinos Patsakis, Quan Le, and Giulio Zizzo. Assessing the impact of packing on static machine learning-based malware detection and classification systems. *Computers & Security*, page 104495, 2025.

[11] Richard Harang and Ethan M Rudd. Sorel-20m: A large scale benchmark dataset for malicious pe detection. *arXiv preprint arXiv:2012.07634*, 2020.

[12] Rakibul Hasan, Barna Biswas, Md Samiun, Mohammad Abu Saleh, Mani Prabha, Jahanara Akter, Fatema Haque Joya, and Masuk Abdullah. Enhancing malware detection with feature selection and scaling techniques using machine learning models. *Scientific Reports*, 15(1):9122, 2025.

[13] Subir Panja, Subhash Mondal, Amitava Nag, Jyoti Prakash Singh, Manob Jyoti Saikia, and Anup Kumar Barman. An efficient malware detection approach based on machine learning feature influence techniques for resource-constrained devices. *IEEE Access*, 2025.

[14] Edward Raff and Charles Nicholas. A survey of machine learning methods and challenges for windows malware classification. *arXiv preprint arXiv:2006.09271*, 2020.

[15] Muhammad Furqan Rafique, Muhammad Ali, Aqsa Saeed Qureshi, Asifullah Khan, and Anwar Majid Mirza. Malware classification using deep learning based feature extraction and wrapper based feature selection technique. *arXiv preprint arXiv:1910.10958*, 2019.

[16] Royi Ronen, Marian Radu, Corina Feuerstein, Elad Yom-Tov, and Mansour Ahmadi. Microsoft malware classification challenge. *arXiv preprint arXiv:1802.10135*, 2018.

[17] Alireza Souri and Rahil Hosseini. A state-of-the-art survey of malware detection approaches using data mining techniques. *Human-centric Computing and Information Sciences*, 8(1):1–22, 2018.

[18] Guillermo Suarez-Tangil, Santanu Kumar Dash, Mansour Ahmadi, Johannes Kinder, Giorgio Giacinto, and Lorenzo Cavallaro. Droidsieve: Fast and accurate classification of obfuscated android malware. In *Proceedings of the seventh ACM on conference on data and application security and privacy*, pages 309–320, 2017.

[19] Nguyen Ngoc Toan, Dang Quang Thang, et al. Static feature selection for iot malware detection. *Journal of Science and Technology on Information security*, 1(15):74–84, 2022.

[20] Yinwei Wu, Meijin Li, Qi Zeng, Tao Yang, Junfeng Wang, Zhiyang Fang, and Luyu Cheng. Droidrl: Feature selection for android malware detection with reinforcement learning. *Computers & Security*, 128:103126, 2023.

[21] Limin Yang, Arridhana Ciptadi, Ihar Laziuk, Ali Ahmadzadeh, and Gang Wang. Bodmas: An open dataset for learning based temporal analysis of pe malware. In *4th Deep Learning and Security Workshop*, 2021.

[22] Muhammed Mutlu Yapici. A novel image based approach for mobile android malware detection and classification. *Knowledge-Based Systems*, page 113855, 2025.

Table 3: Dataset Characteristics and Experimental Configuration

| Characteristic | Microsoft Big-2015 | BODMAS |
|---|---|---|
| *Dataset Properties* | | |
| Training Samples | 8,694 | 94,157 |
| Test Samples | 2,174 | 20,166 |
| Feature Dimensions | 1,795 | 2,381 |
| Classification Type | 9-class | Binary |
| Data Split Ratio | 80:20 | 80:20 |
| *Model Architecture* | | |
| **D3QN (Dueling Double Deep Q-Network)** | | |
| Feature Extraction Layers | Linear(FEATURE_DIM, 128) → PReLU → Linear(128, 128) → PReLU → Linear(128, 128) → PReLU | |
| Value Stream | Linear(128, 1) | |
| Advantage Stream | Linear(128, ACTION_DIM) | |
| Output Layer | Q-values = Value Stream + (Advantage Stream - Advantage.mean()) | |
| **DDQN (Double Deep Q-Network)** | | |
| Hidden Layers | 3 layers, each Linear(128, 128) with PReLU activation | |
| Output Layer | Linear(128, ACTION_DIM) | |
| Architecture Flow | Input → FC Layers → PReLU → Q-value Output | |
| *Training Hyperparameters* | | |
| Training epochs | 10k | |
| Learning Rate (Initial) | $1.0 \times 10^{-3}$ | |
| Optimizer | Adam | |
| Weight Decay | $1.0 \times 10^{-6}$ | |
| Loss Function | Mean Squared Error (MSE) | |
| Gradient Clipping | Max Norm = 1.0 | |
| Target Network Update Rate ($\rho$) | 0.01 | |
| *Learning Rate Scheduling* | | |
| LR Decay Factor | 0.7 | |
| LR Decay Epochs | 3,000 | |
| Minimum LR | $3.0 \times 10^{-8}$ | |
| *RL-Specific Parameters* | | |
| Action Space Dimension | 1,804 (1,795 + 9) | 2,383 (2,381 + 2) |
| State Representation | Feature vector concatenated with selection mask | |
| Episode Termination | Classification decision | |
| Reward Function | Classification accuracy with feature efficiency penalty | |
| *Computational Environment* | | |
| Hardware | CUDA-enabled GPU | |
| Framework | PyTorch | |
| Precision | 32-bit floating point | |

Table 4: Supplementary Material: The following table presents comprehensive performance metrics for all evaluated methods across both datasets, supporting the comparative analysis presented in the main manuscript.

| Method | Big2015 Dataset | | | | BODMAS Dataset | | | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Precision | Recall | F1-Score | Accuracy | Precision | Recall | F1-Score |
| *Traditional Methods (All Features)* | | | | | | | | |
| Decision Tree | 94.25 | 94.09 | 94.25 | 94.11 | 95.76 | 95.76 | 95.76 | 95.76 |
| Random Forest | 98.02 | 97.68 | 98.02 | 97.84 | 98.10 | 98.10 | 98.10 | 98.10 |
| Logistic Regression | 97.65 | 97.96 | 97.65 | 97.74 | 98.19 | 98.21 | 98.19 | 98.20 |
| XGBoost | 98.85 | 98.87 | 98.85 | 98.85 | 98.11 | 98.12 | 98.11 | 98.11 |
| Voting Classifier | 98.16 | 98.11 | 98.11 | 98.11 | 98.50 | 98.50 | 98.50 | 98.50 |
| Stacking Classifier | 97.93 | 98.62 | 97.93 | 98.19 | 98.60 | 98.60 | 98.60 | 98.60 |
| *Feature Selection Methods (60 features)* | | | | | | | | |
| DT + RF Importance | 94.11 | 94.02 | 94.11 | 94.02 | 96.79 | 96.79 | 96.79 | 96.79 |
| RF + RF Importance | 97.88 | 97.56 | 97.88 | 97.71 | 98.00 | 98.01 | 98.00 | 98.00 |
| LR + RF Importance | 94.34 | 95.98 | 94.34 | 94.90 | 93.10 | 93.45 | 93.10 | 93.13 |
| DT + Mutual Information | 91.58 | 91.54 | 91.58 | 91.48 | 96.16 | 96.16 | 96.16 | 96.16 |
| RF + Mutual Information | 95.72 | 95.44 | 95.72 | 95.56 | 97.60 | 97.61 | 97.60 | 97.61 |
| LR + Mutual Information | 86.20 | 88.57 | 86.20 | 86.71 | 85.49 | 86.30 | 85.49 | 85.56 |
| DT + Chi-Squared | 92.18 | 92.31 | 92.18 | 92.08 | 83.69 | 85.06 | 83.69 | 83.78 |
| RF + Chi-Squared | 95.86 | 95.65 | 95.86 | 95.71 | 88.63 | 89.12 | 88.63 | 88.68 |
| LR + Chi-Squared | 92.46 | 93.91 | 92.46 | 92.90 | 89.06 | 90.09 | 89.06 | 89.12 |
| DT + RFE | 95.77 | 95.48 | 95.77 | 95.57 | 97.04 | 97.04 | 97.04 | 97.04 |
| RF + RFE | 98.07 | 97.74 | 98.07 | 97.89 | 97.93 | 97.95 | 97.93 | 97.93 |
| LR + RFE | 94.53 | 95.88 | 94.53 | 94.98 | 94.26 | 94.51 | 94.26 | 94.28 |
| *Proposed RL Methods having episodic length ∼60* | | | | | | | | |
| DDQN (Baseline) | <u>99.12</u> | <u>99.14</u> | <u>99.13</u> | <u>99.13</u> | <u>98.83</u> | <u>98.49</u> | **98.77** | <u>98.63</u> |
| **D3QN (Proposed)** | **99.22** | **99.23** | **99.22** | **99.20** | **98.83** | **98.64** | <u>98.65</u> | **98.64** |

All metrics reported as percentages (%)
**Bold**: Best overall performance per dataset; <u>Underlined</u>: Second-best performing method per dataset
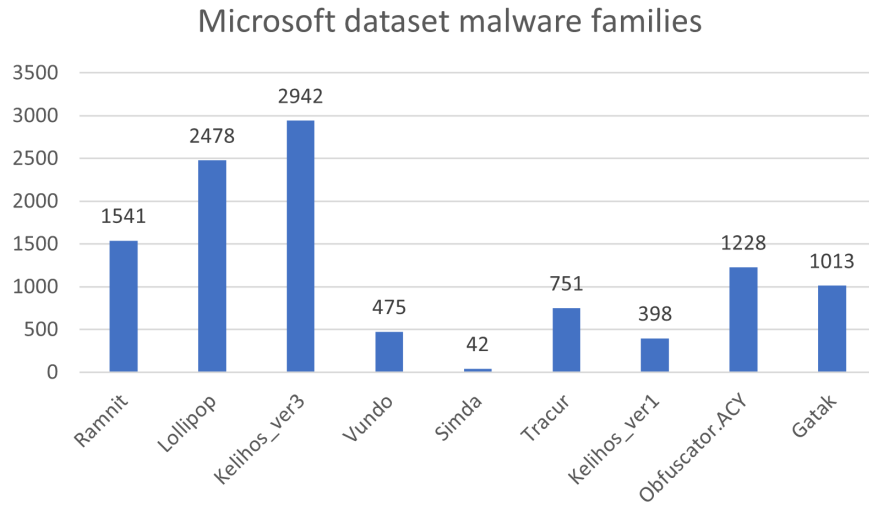DT = Decision Tree, RF = Random Forest, LR = Logistic Regression, RFE = Recursive Feature Elimination

Figure 7: Total number of samples and data distribution of each malware family in the Microsoft Malware Challenge training dataset.
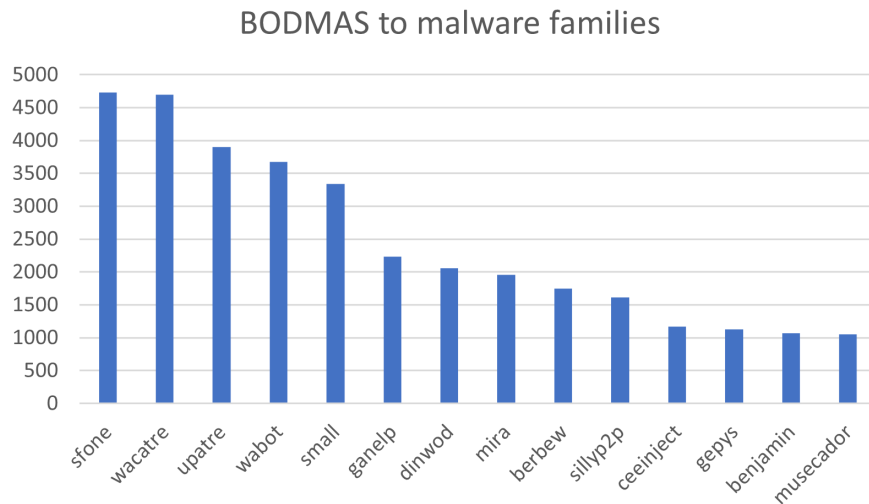


Figure 8: Top 14 malware families and their number of samples ($>= 1,000$) in BODMAS dataset.
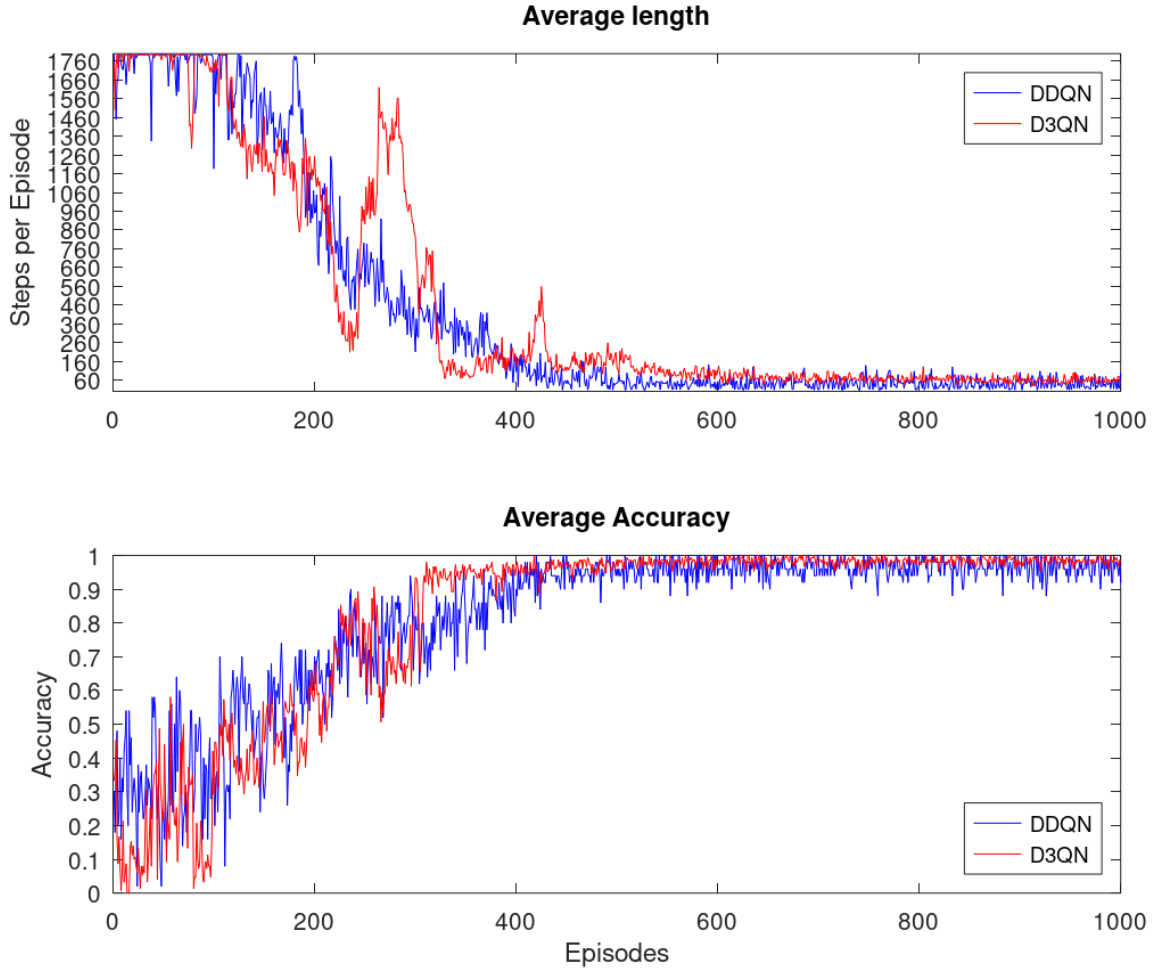
Figure 9: Average Episode Length and Average Accuracy measured on the validation dataset during training for both DDQN and D3QN. These plots illustrate the learning dynamics of the two multi-class classification policy networks trained on the Big2015 dataset using a reduced feature set comprising **1795** features. From the figures, it is evident that our proposed D3QN architecture demonstrates improved performance, achieving higher accuracy with shorter episode lengths. The reduced episode length corresponds to fewer penalties incurred, indicating more efficient decision-making over time.
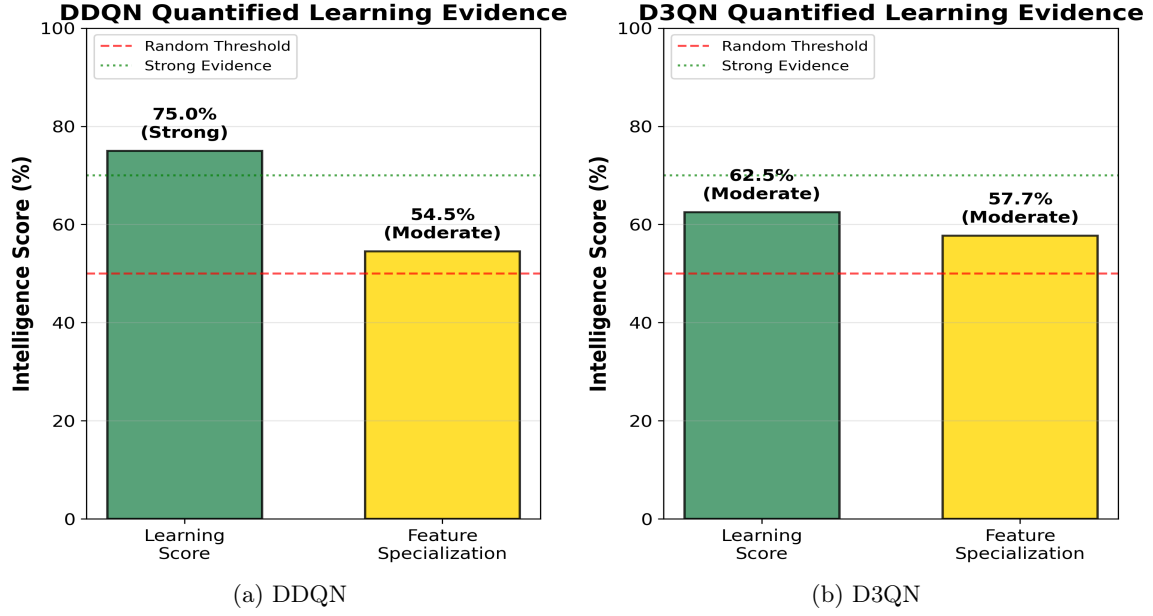
(a) DDQN

(b) D3QN

Figure 10: Quantified Learning Evidence Metrics for both architectures. Learning Score measures the proportion of feature categories showing statistically significant deviations from random baseline (50% threshold), while Feature Specialization quantifies strategic divergence between malware and benign selection patterns. Both DDQN (75.0%, 54.5%) and D3QN (62.5%, 57.7%) demonstrate learned intelligence above random behavior.
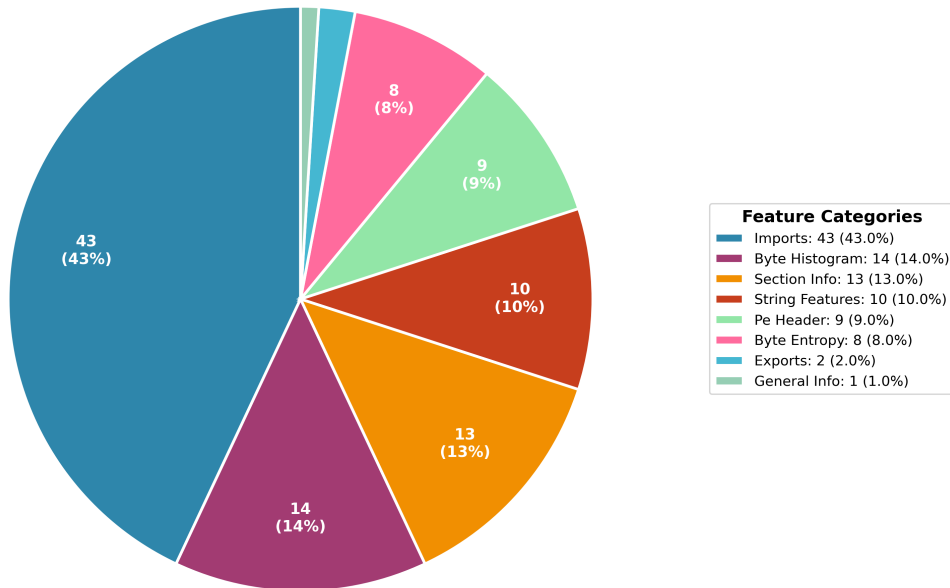


Figure 11: EMBER feature category distribution in top-performing features selected by the model. The chart reveals which feature categories (imports, exports, byte histogram, etc.) are most frequently represented among the highest-importance features, providing insights into the model's learned domain expertise and feature prioritization strategies.
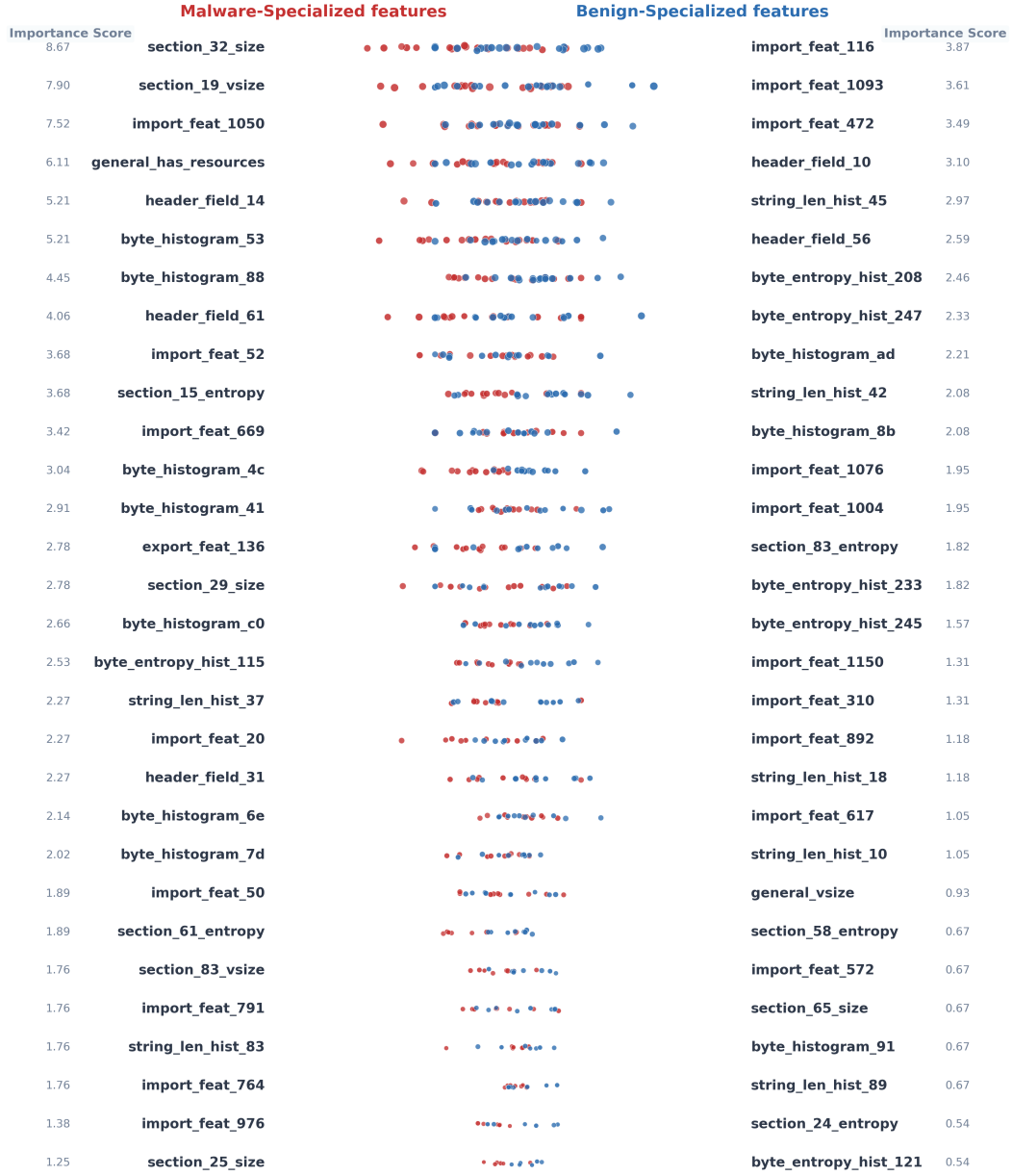
Figure 12: Feature specialization analysis showing discriminative features identified by the reinforcement learning-based malware detection system. Features are ranked by discrimination strength (importance scores shown on sides), with dot intensity and density representing the degree of specialization. **Left (red):** Features preferentially selected by the RL agent for malware classification. **Right (blue):** Features preferentially selected for benign classification. The analysis reveals that malware detection relies heavily on byte histogram and entropy features, while benign classification depends more on import table and header field characteristics. Higher importance scores indicate stronger specialization, with top-ranked features showing the most consistent selection patterns across the episodic decision-making process.