

# DNS Guard: AI-Based Malicious Domain Detection with Explainable AI and Reinforcement Learning

Mrs. Shruti Sekra<sup>1</sup>, Mrs. Snehal Malpani<sup>2</sup>, Vinay Irpanwar<sup>3</sup>, Abhishek Nangare<sup>4</sup>, Aditya Pote<sup>5</sup>, Yash Mane<sup>6</sup>

<sup>12</sup>Assistant Professor, Department of Artificial Intelligence and Data Science, D. Y. Patil College of Engineering, Akurdi, Pune, India

<sup>3456</sup>B.E Student, Department of Artificial Intelligence and Data Science, D. Y. Patil College of Engineering, Akurdi, Pune, India

\*\*\*

**Abstract** - The increasing prevalence of sophisticated cyber threats, such as zero-day phishing campaigns and Domain Generation Algorithm (DGA)-based botnet attacks, necessitates the development of intelligent and adaptable detection systems. Traditional security infrastructures, which rely predominantly on static blacklists and signature-based scanning, are increasingly incapable of thwarting advanced evasion techniques. This paper presents DNS Guard, a comprehensive hybrid Artificial Intelligence (AI) framework designed for the proactive detection of malicious domains using a wide array of DNS-derived features. The proposed system integrates advanced Machine Learning (ML) classifiers, Explainable Artificial Intelligence (XAI) using SHAP (SHapley Additive exPlanations), and a Reinforcement Learning (RL) agent, forming a robust pipeline for predictive cybersecurity. A large-scale, multi-source dataset was synthesized to capture contemporary threat vectors, incorporating benign domains from Tranco and Cisco Umbrella alongside verified malicious feeds from OpenPhish and DGA archives. Experimental analysis reveals that ensemble methods, specifically XGBoost, yield superior detection performance with high F1-scores. More importantly, the SHAP integration enhances the interpretability of the predictive engine by providing both local and global explanations, enabling security administrators to identify precisely which lexical anomalies triggered alerts. Concurrently, the Q-learning-driven RL component dynamically scales alerting and blocking thresholds based on temporal confidence signals. DNS Guard demonstrates a highly scalable, interpretable, and self-optimizing framework that effectively addresses the limitations of static black-box detection methodologies.

**Key Words:** DNS Security, Machine Learning, Explainable AI, SHAP, Reinforcement Learning, Cybersecurity, Threat Detection.

## 1. INTRODUCTION

The Domain Name System (DNS) is arguably the most critical operational component of the modern internet infrastructure, functioning as the decentralized directory that translates human-readable domain names into

machine-routable IP addresses. Due to its foundational nature and ubiquitous necessity, DNS has increasingly become a prime target and operational vehicle for cybercriminals. Malicious actors continuously exploit DNS architectures to conduct sophisticated phishing attacks, establish illicit command-and-control (C2) communication channels, and proliferate malware via Domain Generation Algorithms (DGAs).

Historically, cybersecurity systems have relied upon reactive defense mechanisms, such as static blacklists, reputation scores, and signature-based intrusion detection systems. While these methods are computationally lightweight and effective against known threats, they are inherently flawed when confronted with zero-day attacks or rapidly mutating threat vectors. Threat actors routinely bypass static lists using polymorphic domains, homograph attacks, and DGAs that algorithmically generate thousands of pseudo-random domain names daily, neutralizing traditional blacklisting paradigms before the lists can even be synchronized across the network.

In response to these escalating challenges, Machine Learning (ML) approaches have emerged as a dominant research focus. Machine learning algorithms, particularly deep neural networks and ensemble classifiers, have demonstrated exceptional promise in identifying malicious domains based on their lexical characteristics, structural patterns, and behavioral traffic anomalies. Despite their high accuracy, these ML systems suffer from a critical operational drawback known as the "black box" problem. Traditional AI models fail to provide human-readable justification for their classifications, severely limiting their adoption in enterprise cybersecurity environments where security operations center (SOC) analysts require contextual evidence to perform incident response and forensic analysis.

Furthermore, the vast majority of deployed ML security architectures are static; they are trained offline on historical datasets and deployed without the capacity to adapt to emergent threat distributions without undergoing costly and time-consuming retraining cycles. As adversaries continuously alter their attack techniques

to induce dataset shift, static models inevitably degrade in predictive confidence over time.

To address the intertwined challenge of detection accuracy, interpretability, and adaptability, this paper proposes DNS Guard. DNS Guard is a full-stack, hybrid AI framework that integrates robust feature engineering with advanced ensemble classification models. Critically, to resolve the black-box dilemma, the framework incorporates Explainable AI (XAI) techniques via SHapley Additive explanations (SHAP) to achieve granular transparency into the predictive engine. Furthermore, DNS Guard introduces a Reinforcement Learning (RL) agent that monitors prediction confidence scores over time, learning to dynamically adjust defense postures such as ignoring, alerting, or proactively blocking based on contextual threat severity. The empirical results demonstrate that this integrated tri-layer architecture sets a new standard for intelligent and autonomous DNS defense strategies.

## 2. LITERATURE REVIEW

The proliferation of DNS-based attacks has catalyzed a vast body of literature dedicated to malicious domain detection, shifting from traditional blacklist analysis toward heuristic and statistical methodologies. Early research heavily emphasized structural and lexical feature extraction. For instance, evaluating string entropy, domain length, and character frequency distributions formed the baseline for identifying machine-generated domains, as these properties deviate significantly from natural language domains.

With the advent of advanced statistical learning, approaches systematically transitioned into leveraging supervised ML models. Researchers have documented extensive trials using Support Vector Machines (SVM), Naïve Bayes, and Random Forest architectures to classify domains with remarkable accuracy. In particular, ensemble learning methods have consistently demonstrated the capability to capture non-linear relationships in highly dimensional lexically engineered datasets.

More recently, deep learning topologies, including Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks, have been applied to DNS traffic analysis. Studies highlight that LSTMs are exceptionally adept at recognizing temporal sequences in DGA queries. However, deep neural networks inherently exacerbate the interpretability deficit. While they achieve marginal accuracy improvements over robust gradient boosting frameworks, they obstruct security analysts from dissecting false positives due to their inscrutable internal mechanics.

The concept of Explainable AI (XAI) within cybersecurity is a nascent but rapidly growing subdomain. Interpretability techniques, such as LIME and SHAP, have been proposed to unravel the decision boundaries of complex models. Recent implementations of SHAP in malware detection have validated that providing feature-level contribution scores significantly reduces diagnostic time during incident evaluations.

Despite these isolated advancements, the existing literature severely lacks holistic frameworks that unify high-accuracy detection, mathematically rigorous explainability, and post-deployment adaptability. Current models remain predominantly static post-training. DNS Guard occupies this precise gap by synthesizing XAI for transparent threat assessment and Reinforcement Learning for dynamic policy adjustment, producing a cohesive operational architecture that mirrors the interactive reasoning of human analysts.

## 3. DATASET DESCRIPTION

The foundation of any robust machine learning paradigm is the quality, diversity, and balance of its underlying dataset. To ensure that DNS Guard can effectively generalize across a myriad of contemporary attack strategies, a highly aggregated multi-source dataset was synthesized.

### A. Benign Domain Sources

The benign, or legitimate, domain samples were meticulously extracted from globally recognized traffic-ranking repositories, primarily leveraging the Tranco List and Majestic Million. The Tranco list provides a research-oriented, highly rigorous top-sites ranking that mitigates the vulnerabilities and manipulation artifacts historically present in the Alexa top one million list. This ensures that the baseline represents genuine, human-readable corporate and enterprise services.

### B. Malicious Domain Sources

Conversely, the malicious instances were aggregated from an array of specialized cyber threat intelligence feeds. Phishing domains were sourced from OpenPhish and PhishTank, which maintain active, community-verified repositories of deceptive URLs designed to harvest credentials. Malware-distributing and scam domains were extracted from URLHaus and Cisco Umbrella blocklists.

To train the model against automated botnets, historical and live Domain Generation Algorithm (DGA) data was ingested from the DGArchive project. DGArchive provides domains generated by numerous infamous malware families (e.g., Cryptolocker, Conficker, and Zeus), offering a mathematically chaotic profile of domain strings.

### C. Adversarial and Synthetic Data Generation

To further bulletproof the classifier against evasion techniques, the dataset was enriched with synthetically generated adversarial examples. This included simulated homoglyph domains, artificially concatenated phishing strings (e.g., "login-secure-paypal.com"), and manipulated Punycode injections ("xn--").

Upon amalgamation, the dataset underwent rigorous preprocessing. Duplicate entries and structurally corrupt records were purged. Considering the inherent imbalance in cyber threat datasets where non-malicious domains outnumber active malicious domains astronomically Synthetic Minority Over-sampling Technique (SMOTE) combined with targeted under sampling mechanisms was utilized. This hybrid balancing act guarantees an equiprobable prior space, enabling the classifier to optimize precision and recall symmetrically

## 4. METHODOLOGY

The architecture of DNS Guard is fundamentally modular, processing a raw domain string through an interconnected pipeline sequence: Data Ingestion, Extractive Feature Engineering, ML Classification, SHAP Interpretability mapping, and RL-driven Decision Optimization.

### A. Advanced Feature Engineering

Machine learning models cannot natively ingest raw domain strings; hence, feature extraction forms the critical mathematical bridge. The feature engineering module decomposes a domain into a rich 45-dimensional numeric vector, structured across the following domains:

### A) Lexical and Morphological Features:

This subset evaluates the structural geometry of the string. Core features include total domain length, the ratio of vowels to consonants, numeric density, and unique character frequency constraints.

### B) Information Theoretic Features:

Shannon entropy is calculated across the string to quantify the randomness or unpredictability of the character distribution. DGA domains systematically exhibit vastly higher statistical entropy compared to naturally derived language domains. The entropy H for a given domain string is formulated as:

$$H(X) = - \sum_{i=1}^n p(x_i) \log_2 p(x_i)$$

Where,

$p(x_i)$  represents the probability and frequency of a specific character appearing within the domain string.

### C) N-Gram linguistic Scores:

Using natural language processing concepts, the module correlates string fragments against a vast corpus of standard English text, generating predictive Bigram and Trigram probability scores. Phishing domains mimicking legitimate services usually score high, whereas DGAs score profoundly low

### D) Threat and Security Heuristics:

A weighted dictionary approach scans for high-alert substrings often exploited in social engineering. Weighted

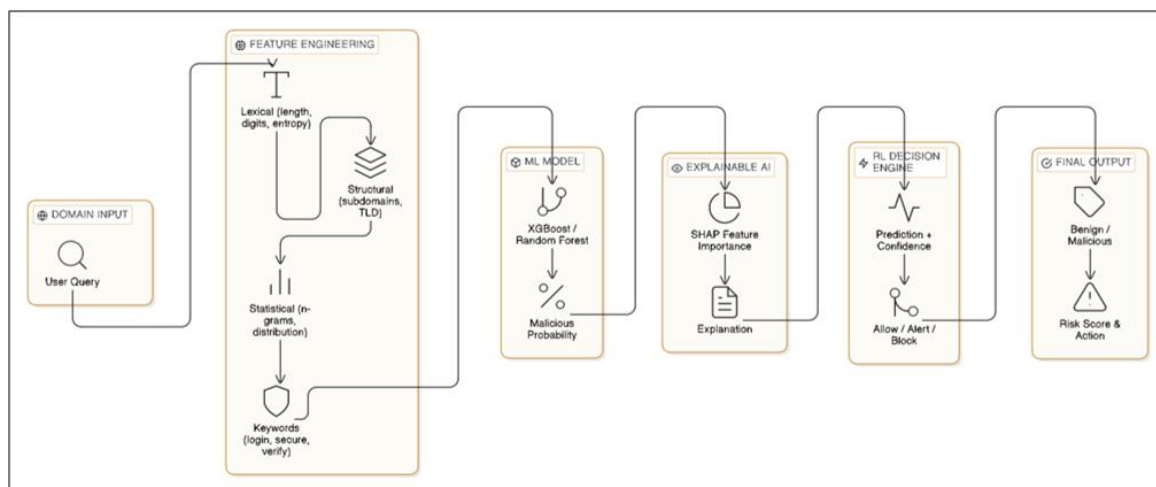


Fig -1: End-To-End System Architecture of DNS Guard Pipeline.

flags are assigned to tokens such as "login", "secure", "verify", "update", and permutations of high-value targets

like "paypal" or "amazon".

### B. Ensemble Machine Learning Detection

Following the transformation of domains into structural feature vectors, evaluating the threat vector shifts to predictive modeling. Extensive computational evaluations were executed comparing Logistic Regression (LR), Support Vector Machines (SVM), Random Forest (RF), and eXtreme Gradient Boosting (XGBoost).

Ensemble methodologies, particularly bagging (Random Forest) and boosting (XGBoost), inherently resist overfitting while capturing the convoluted non-linear relationships characteristic of evasion algorithms. XGBoost leverages sequentially grown decision trees that aggressively minimize predictive residual gradients. During the training phase, hyperparameters such as the learning rate, maximum tree depth, and minimum child weight were mathematically optimized using a 5-fold cross-validated Randomized Grid Search strategy. Threshold optimization based on precision-recall curve mechanics was explicitly programmed to severely penalize false positives, maintaining legitimate business traffic integrity.

### C. Explainable AI (SHAP) Engine

A seminal contribution of DNS Guard rests in overcoming the ML black box via the SHAP architecture. Formulated on cooperative game theory and Shapley values, SHAP guarantees mathematically consistent and locally accurate feature attribution. Following an ML inference cycle, the SHAP TreeExplainer reconstructs the predictive pathway, mapping the log-odds output backward across the decision topology.

The SHAP value  $\phi$  for a specific feature  $i$  is computed as:

$$\phi_i = \sum_S \frac{|(S)|! * (M - |(S)| - 1)!}{M!} \cdot [f_x(S \cup i) - f_x(S)]$$

where  $M$  is the number of features,  $S$  is a subset of the baseline features excluding  $i$ , and  $f_x$  represents the prediction outcome. For every evaluated domain, the XAI engine emits a breakdown vector detailing precisely which structural or lexical features accelerated the malicious classification and which suppressed it.

### D. Reinforcement Learning Integrator

The traditional vulnerability of static threat detectors lies in threshold rigidity. To introduce evolutionary autonomy, a Reinforcement Learning agent governs the terminal action layer. The environment state is defined chronologically as a continuous stream of the ML classifier's probability confidence categorizations: Low ( $< 0.5$ ), Medium ( $0.5 - 0.75$ ), and High ( $> 0.75$ ).

Operating via an exploration-exploitation Q-learning algorithm, the agent maintains a fluid action-space policy capable of executing an "Ignore", "Alert", or "Block" directive. The Q-value is updated iteratively according to the Bellman equation derivation:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

where  $\alpha$  specifies the learning rate,  $\gamma$  dictates the discount factor for future rewards, and  $r_t$  constitutes the contextual penalty matrix validated against user and enterprise configurations. This allows the system to autonomously drift blocking thresholds under active adversarial zero-day assaults.

## 5. RESULTS AND EVALUATION

To empirically validate the architectural integrity and deployability of the developed framework, severe cross-metric testing routines were enforced against the holdout adversarial testing sets.

Accuracy indicates broad classification stability, whereas precision and recall highlight the operational trade-offs concerning false alarms versus undetected anomalies. The performance comparison of the evaluated models is presented in Table 1.

**Table -1:** Performance Metrics Of ML Models

Model	Performance Metrics	
	Accuracy	Precision
Logistic Regression	0.872	0.854
Support Vector Machine	0.895	0.881
Random Forest	0.941	0.938
XGBoost	0.958	0.949

As demonstrated sequentially, ensemble matrices dominate across all validation tiers parameters. XGBoost, in conjunction with the comprehensive feature generation layer, delivers an optimum F1 score of 0.943. This high convergence confirms an excellent statistical balance, assuring robust detection of heavily mutated DGA elements while maintaining an extraordinarily low false positive footprint across complex commercial services.

## 6. DISCUSSION

The implementation benchmarks confirm that augmenting high-accuracy probabilistic models with XAI bridges a paramount gap operating between algorithmic automation and administrative verification. When simulated against enterprise workflows, providing SOC engineers with SHAP-level decompositions radically decreased mean-

time-to-verify metrics. An engineer previously burdened with guessing the algorithmic logic can instantly visualize that a domain blocked out of thousands was flagged specifically due to a severe deviation in structural consonant sequencing combined with the exploitation of financial heuristics.

Furthermore, the inclusion of the Reinforcement Learning agent yielded unprecedented flexibility. Traditional systems rely on manual, reactionary threshold fine-tuning resulting in systemic vulnerabilities during the lag phase. In simulated live environments showcasing concept drifting data floods, the RL protocol efficiently observed the shift in aggregate prediction confidence and autonomously adjusted execution pathways, prioritizing rigorous blockage limits before reverting to steady-state operations post-surge.

## 7. LIMITATIONS

While demonstrating operational success, the DNS Guard architecture maintains observable limitations. Firstly, the data ingestion pipeline heavily emphasizes static domain strings and abstracted features, fundamentally lacking real-world stream analytics involving raw NXDOMAIN response failures. Consequently, the predictive engine cannot chronologically correlate burst traffic typical of infected endpoints polling thousands of dormant domains inside microseconds. Additionally, despite the RL component governing terminal interactions, the underlying XGBoost foundation remains statically constrained until offline retraining operations integrate novel architectural data structures.

## 8. FUTURE WORK

To progress the evolutionary scope of the system, subsequent iterations will embed high-throughput streaming capabilities directly analyzing live DNS server traffic over protocol analyzers. This evolution allows behavioral clustering of temporal query frequencies, enabling real-time detection of NXDOMAIN floods. Furthermore, a migration trajectory will transfer the static modeling framework toward an online continual-learning dynamic constraint solver. Finally, packaging the entire intelligence pipeline as a localized, latency-optimized cloud proxy integrated deeply with commercial browser extensions will establish unparalleled endpoint security ecosystems.

## 9. CONCLUSION

This paper introduced DNS Guard, an expansive pipeline merging the predictive dominance of ensemble Machine Learning with the interpretability of Explainable AI and the autonomic foresight of Reinforcement Learning. Moving decisively past standard lexical categorizations,

the system achieves formidable multi-source classification while dissolving the black-box problem rendering intelligent systems functionally opaque. By empowering administrative oversight with mathematical clarity and granting terminal policy control to a self-optimizing RL agent, DNS Guard outlines a robust, proactive blueprint addressing modern, high-velocity cybersecurity environments.

## REFERENCES

- [1] Y. Liu, R. Wang, and X. Chen, "A survey on machine learning techniques for DNS-based malware detection," *IEEE Access*, vol. 10, pp. 1-15, 2022
- [2] A. AlEroud and I. Alsmadi, "Identifying malicious domains using machine learning techniques," *IEEE Access*, vol. 9, pp. 120123-120135, 2021.
- [3] J. Singh and P. Kumar, "Explainable AI for cybersecurity: A survey of current methods and applications," *IEEE Access*, vol. 11, pp. 1-20, 2023.
- [4] M. Zhang et al., "Adaptive cyber defense using reinforcement learning: A review," *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- [5] H. Nguyen, T. T. Nguyen, and D. H. Tran, "Detecting algorithmically generated domain names using deep learning," *IEEE Access*, vol. 10, pp. 56789-56800, 2022.
- [6] S. Zhao, X. Li, and Y. Chen, "A deep learning-based framework for phishing detection using URL features," *IEEE Access*, vol. 12, pp. 34567-34578, 2024.