# Countering DNS Vulnerability to Attacks Using Ensemble Learning

Love Allen Chijioke Ahakonye, Cosmas Ifeanyi Nwakanma, Simeon Okechukwu Ajakwe*,
Jae Min Lee, and Dong-Seong Kim
*IT Convergence Engineering*, *Kumoh National Institute of Technology* Gumi, South Korea
(loveahakonye, cosmas.ifeanyi, ljmpaul, dskim) @kumoh.ac.kr
simeonajlove@gmail.com*

*Abstract*—The Domain Name System (DNS) is the hub of the cyberspace and communications services which also plays enabling role in the Industrial Internet of Things (IIoT) and transmission at large. DNS enciphering in HyperText Transfer Protocol Secure (HTTPS) as DoH did not eliminate vulnerability and intrusion into critical systems. This study proposed a time-efficient Ensemble Learning (EL) model for countering DNS vulnerability to attack. The proposed EL candidate incorporates feature selection capability in extracting relevant features for enhanced model optimization. The simulation results showed that the proposed EL candidate effectively mitigates vulnerability, classifying DNS traffic into Non-DoH, Malicious DoH and Benign-DoH. The proposed model outperforms other compared state-of-art EL techniques with a combined advantage of accuracy and training time of $99.5\%$ and $13.96s$.

*Index Terms*—AI, DNS, Ensemble Learning, IIoT,

## I. INTRODUCTION

Amongst the protocols in a network communication system is the Domain Name System (DNS). It serves as the internet directory. Online information access is through the DNS. It is also one of the early network protocols with high vulnerability and diverse security flaws constantly exploited. DNS exploitation is invariably a domain of eminent attention for cyber-security researchers. However, delivering privacy and safeguarding DNS demands and acknowledgement remains a daunting endeavor as intruders employ advanced intrusion strategies for exploiting DNS vulnerability [1].

Some of the DNS attacks are domain lock-up attacks, DNS hijacking, DNS spoofing, DNS Tunneling. In pursuit of security improvement, DNS has become more relevant recently by providing validation and approval to some internet resources. Nevertheless, the recent development of DNS falls short of guaranteeing requisite security to users. Hence, DNS security has been a widely researched topic in the cyber-security domain. The National Institute of Standards and Technology (NIST) issued a document with recommendations for the safe deployment of DNS to prevent security challenges [2].

To reduce the DNS vulnerabilities associated with security and data processing, the Internet Engineering Task Force (IETF) initiated DNS over HTTPS (DoH) in RFC8484 [1]. DoH is a set of codes that strengthens security and counters attacks by encoding DNS mistrust and forwards in a hidden channel such that there is no data obstruction in transit. How-

ever, the inadequacy of an illustrative dataset is a challenge in evaluating the approach for securing DoH traffic in a network topology. Distinct perspectives for DNS vulnerability such as Internet Protocol (IP)/domain boycott and removing suspected DNS packets to attain DNS restriction have been applied [3]–[5]. Most researchers criticize DoH for making DNS tunnels difficult for attack detection.

There have been efforts at protecting network systems; this attempt includes network intrusion detection systems (NIDS), the use of firewalls to minimize the issues of unlawful access, etc. Consequently, the authors in [1] attempted resolving the problem of dataset insufficiency. The study proposed an approach to secure a model dataset. It is for the analysis, testing and evaluation of DoH traffic in hidden channels. The focus was on deploying DoH in an application to take hold of benign and malicious DoH traffic. This new protocol improved privacy and DNS security. However, the issue of DNS vulnerability persists, hence, the need for an Artificial Intelligence (AI) countering measure.

AI has supported in strengthening the performance of detection techniques in NIDS. Considering this, numerous authors have utilized AI techniques for vulnerability and attack mitigation. For instance, authors [6]–[8] proposed various intrusion detection frameworks. Presently industrial innovations such as the industrial internet of things (IIoT) and machine learning (ML) have revolutionized daily life and influenced various sectors. IIoT has become ubiquitous as it is applicable in different areas, including communication and the industry generally. ML has played enabling role in the development of Intrusion Detection Systems (IDS). Such application includes attack and vulnerability detection, which has found use in IIoT.

The use of ML for enabling attack and vulnerability detection is still a contending issue. Thus, this study investigated the DoH traffic and non-DoH traffics for an efficient IDS. Also, an investigation into different ensemble learning (EL) prospects for an efficient, accurate and time-aware countering system for DNS vulnerability to attacks. Recent research works show that ML would not only improve the detection rate but would also reduce the computational time [9].

This work has the following goals:
- To deduce an efficient AI technique for IDS in terms of the combined advantage of accuracy and least training time using the CIRA-CIC-DoHBrw-2020 datasets.

- To utilize Python to determine the best EL candidate with respect to DNS vulnerability to attacks.
- To propose an EL architecture to counter DNS vulnerability, achieve high accuracy of detection with reduced complexity and less training time.

The paper arrangement is thus: Section II is the summary of existing works detailing IIoT Intrusion detection and various approaches of EL and establishing research gaps. In Section III, problem formulation, which involves a brief description of EL and the best AI candidate, has been discussed. Section IV describes the performance evaluation with the comparison of accuracy plot and time-efficiency of the proposed model, Section V is the conclusion of the paper.

## II. RELATED WORKS

### A. Intrusion Detection System applying Ensemble Learning (EL)

Indications are replete in literature to deduce that conventional ML approaches may not be reliable in manipulating intricate data, such as high-dimensional data, noisy (usually from industrial environment) and imbalanced data. To solve this challenge, various works on EL schemes are now accessible to ensure a blend of data mining,fusion and modeling into a consolidated scheme [10].

EL is commonly a collective classifier system. It requires the combination and training of collective learners known as base learners, to determine a learning problem. The EL framework can be homologous or diverse ensembles reliant on the type of base learners constituting the scheme. While the homologous have base learners, the diverse ensemble comprises of individual unique learners or simply called component learners. Significantly, most studies on EL schemes are focused on weak learners, thus, base learners are often referred to as weak learners. Fig. 1 depicts the fusion of base learners for an improved outcome, where the outcome is anticipated to be of an enhanced performance in comparison with base learners 1, 2, 3 ...K. The base learners can be perceived as those learners



Fig. 1. Flow of Ensemble Learning showing the fusion of base learners

whose accuracy in terms of binary classification ability is marginally within 50%, [10].

EL techniques uses a blend of distinct classifiers for detection, which has facilitated different operations and improved performance in IoT systems. This yields enhanced performance for various attack types and protocols used in IoT networks. The study by [11] proposed a modern ensemble IDS approach for attack defence on Ethernet Consist Networks of trains. Though this system delivered superior results, the approach can be complex with a lack of computational speed. Authors [12] proposed an AdaBoost EL scheme for mitigating malicious activities, especially HTTP, MQTT and DNS attacks from botnet attacks in IoT networks. The approach holds good performance accuracy when compared to other models. However, the computation time is not efficient for a time-critical system. Recently, the works of [13] proposed a novel scheme named ElStream for detecting concept drift utilizing traditional ML approach and EL. This approach uses only optimum classifier to vote for decision by using the majority voting scheme. Authors [14] in a recent study attempted AI techniques for mitigating attacks in SCADA Systems.

IIoT are real-time systems which are time-critical and as such a vital factor in its design. Regardless of the value improvement of these authors, the techniques lacks consideration for EL technique for mitigating DNS vulnerability and time complexity for proposed EL in attack detection. Moreover, all enumerated studies did not establish an efficient and time-aware EL. Hence, this study presents the application of an effective EL approach for countering DNS attacks.

## III. METHODOLOGY

### A. DNS Traffic Dataset

This dataset is made of recent attack features as enumerated in [1]. Below is a brief overview of the composition of the dataset. Benign-DoH: This is a non-malicious DoH activities gathered from websites that uses HTTPS, and tagged Benign-DoH. In an effort generate sufficient data to stabilize the dataset, numerous webpages from Alexa were surfed. Non-DoH: This data was created with similar method as in Benign-DoH by employing browsers such as Google Chrome and Mozilla Firefox. Malicious-DoH: DNS channeling mechanism like DNSCat2, Iodine and dns2tcp were used to create malicious DoH data. This mechanism sends encoded TCP data in DNS query. Particularly, the mechanism generate channels of encoded traffic. Consequently, DNS query is made by applying traffic layer-encoded HTTPS application to dedicated DoH servers.

### B. Ensemble Learning Framework For DNS

An EL approach is presented to detect vulnerability which reveal IIoT networks through the DNS codes by examining the DNS codes, as displayed in Fig 2. The architecture comprises of three stages, viz: a feature lay, feature selection, and EL techniques. The feature lay is simply the features of the DNS traffic dataset. Following is the Pearson Coefficient Correlation (PCC). It is used for choosing the lowest correlated features with likely features of benign, Non-DoH and malicious patterns as seen in equation 1, with the option of variables that has a high correlation value threshold of between +/-1. Lastly, the EL technique used for countering the vulnerability in
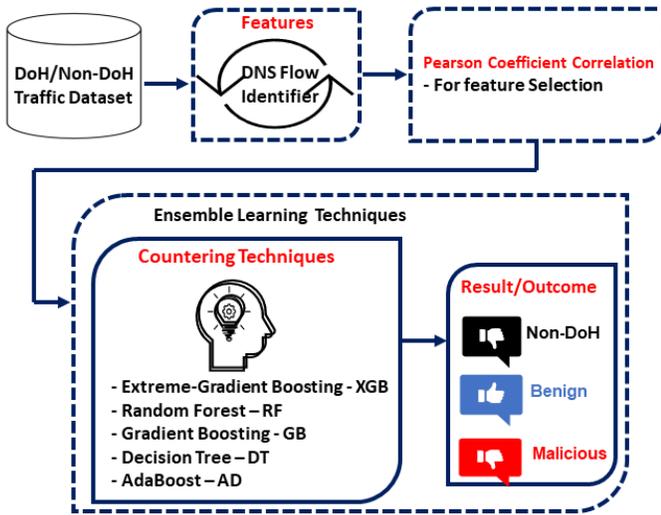
Fig. 2. Proposed Framework for Countering DNS Vulnerability

DNS thereby classifying into Non-DoH, benign and malicious. The dataset was split into training and testing to conduct the experimental evaluation, where 80% of data is used for training and 20% for testing.

$$Q = \frac{\sum \left(\alpha_i - \hat{\alpha}\right)\left(\beta_i - \hat{\beta}\right)}{\sqrt{\sum \left(\alpha_i - \hat{\alpha}\right)^2 \left(\beta_i - \hat{\beta}\right)^2}}, \tag{1}$$

Feature selection is vital in IDS for ridding unnecessary features and choosing the relevant ones that supports segregation of DNS traffic into benign, Non-DoH and malicious. It improves the general performance of the system, lowering computational cost, eliminating information redundancy and enhances accuracy and also helps in the analysis of network data normality. The training specifications are as presented in Table I. This study focused on state-of-the-art EL candidates such as Extreme-Gradient Boosting (XGB), Gradient Boosting (GB), AdaBoost (AD), Random Forest (RF) and Decision Trees (DT) EL classifiers.

TABLE I
ENSEMBLE LEARNING TRAINING PARAMETERS

| Parameter | Remark |
|---|---|
| Observations | 226406 samples |
| Predictors | 11 |
| Classes | 3 |
| No of Trainable Classifiers | 5 |
| Model type | Decision Tree Classifier |
| Result Presentation type | Response plot |
| Training time | 13.960 sec |
| No of Splits | 10 |
| Random State | 42 |
| Cross-validation | kfold |

## IV. PERFORMANCE EVALUATION

### A. Parameter Metrics

To determine an efficient EL technique for countering DNS vulnerability to attacks, the CIRA-CIC-DoHBrw-2020 dataset was evaluated utilizing XGB, GB, AD, RF and DT EL candidates. See Figs 3 and 4 for the performance comparison of the evaluated models. The performance of the prospective EL model was compared with the studies by [11]–[13] using performance evaluation metrics summarized by equations (2) and computation time. The authors in [13] achieved a higher accuracy. However they did not consider computational time as an important factor for such a system. Thus, a research gap which necessitated this new approach.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}, \tag{2}$$

where $FP$, $TP$, $TN$ and $FN$ represents False Positive, True Positive, True Negative and False Negative respectively.
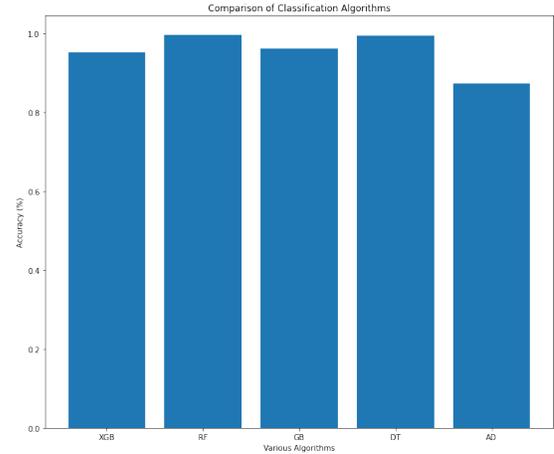


Fig. 3. Comparison of Accuracy Performance of Various EL Techniques
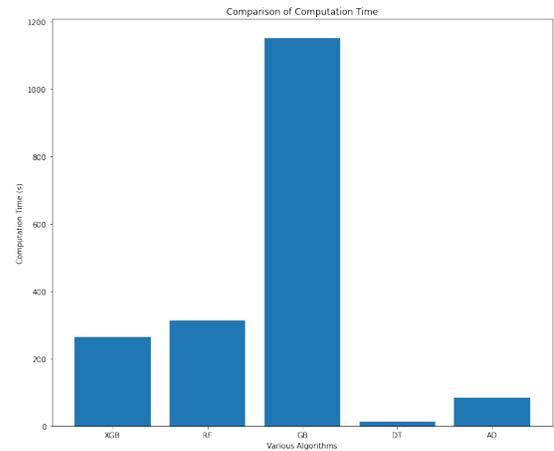


Fig. 4. Comparison of Computation Time of Various EL Techniques

## B. Experimental Environment

Training and testing of the proposed EL scheme was carried out on Google Colab using various Keras and scikit learn libraries. All experimentation can be administered on a laptop designed by NVIDIA GeForce GTx 1080Ti, 11G memory, Intel Xeon E5-1650 CPU 3.60-GHz processor, with windows 10 64 bit operating system.

## C. Comparison of the Performance of some EL Techniques

This study centres on the choice of an efficient EL technique for countering DNS vulnerability to attacks, with effort at comparing the performance of various EL techniques in review literature. Highlighting the achievements of the utilization of EL techniques such as AdaBoost EL (AD-EL), ensemble IDS (E-IDS) and ElStream. Table II, gives an illustration of the achievements based on computation time, accuracy, precision and recall.

TABLE II
PERFORMANCE COMPARISON OF VARIOUS ENSEMBLE TECHNIQUES

| Models | Acc. (%) | Prec (%) | Recall (%) | Comp. Time (ms) |
|---|---|---|---|---|
| DT | 99.3 | 99.2 | 99.3 | 13.96 |
| AD | 87.2 | 87.5 | 86.9 | 84.38 |
| AD-EL [12] | 99.54 | 98.62 | 98.93 | 150.8 |
| XGB | 95.1 | 95.7 | 95.1 | 263.36 |
| GB | 96.0 | 96.0 | 95.6 | 1149.92 |
| RF | 99.5 | 99.4 | 99.6 | 313.23 |
| Elstream [13] | 99.99 | 99.95 | 99.97 | - |
| E-IDS [11] | 97 | 96.8 | 97.5 | - |

## D. Trade-off of Computation Time and Accuracy

It is vital to note that the efficient performance of a model is not solely based on accuracy, rather on a combination of performance evaluation metrics as portrayed in this study. Time is an important factor for time-critical system as DNS security, hence requires the mitigation technique to act as swift as possible. Since any time lapse in mitigating a security breach could lead to fatality in exposure/loss of vital information. On this note, a trade-off between accuracy and time is used as performance metrics. The DT compensated its accuracy (99.3%) with computation time (13.96ms), while other models (AD-EL 99.54%) had longer computation time.

## V. CONCLUSION

This work evaluated various EL candidates leveraging DNS traffic data. The result shows that the performance of the proposed EL DT exceeded other states of the art EL candidates such as RF, XGB, GB and AD for the efficient mitigation of DNS vulnerability, as can be seen in the least computation time of $13.96ms$ and accuracy of $99.3\%$. The specific, practical significance of this ascertainment is in the decision of an efficient EL candidate for IDS, basically where the preference is time-efficiency. For future directions, expanding the comparison by unveiling the capability and flexibility of the DT parameters to more current cyber-security datasets looks promising.

## REFERENCES

[1] M. MontazeriShatoori, L. Davidson, G. Kaur, and A. H. Lashkari, "Detection of DoH Tunnels using Time-series Classification of Encrypted Traffic," in *The 5th IEEE Cyber Science and Technology Congress, Calgary*, 06 2020.

[2] "Secure Domain Name System (DNS) Deployment Guide, author=Chandramouli, Ramaswamy and Rose, Scott," *NIST Special Publication*, vol. 800, pp. 81–2, 2006.

[3] M. MontazeriShatoori, L. Davidson, G. Kaur, and A. Habibi Lashkari, "Detection of DoH Tunnels using Time-series Classification of Encrypted Traffic," in *2020 IEEE Intl Conf on Dependable, Autonomic and Secure Computing, Intl Conf on Pervasive Intelligence and Computing, Intl Conf on Cloud and Big Data Computing, Intl Conf on Cyber Science and Technology Congress (DASC/PiCom/CBDCom/CyberSciTech)*, 2020, pp. 63–70.

[4] A. Nadler, A. Aminov, and A. Shabtai, "Detection of Malicious and Low Throughput Data Exfiltration over the DNS Protocol," *Computers and Security*, vol. 80, pp. 36–53, 2019.

[5] C. Patsakis, F. Casino, and V. Katos, "Encrypted and Covert DNS Queries for Botnets: Challenges and Countermeasures," *Computers and Security*, vol. 88, p. 101614, 2020. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S016740481831321X

[6] G. C. Amaizu, C. I. Nwakanma, S. Bhardwaj, J. M. Lee, and D. S. Kim, "Composite and Efficient DDoS Attack Detection Framework for B5G Networks," *Computer Networks*, vol. 188, p. 107871, 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1389128621000438

[7] M. Teixeira, T. Salman, M. Zolanvari, R. Jain, N. Meskin, and M. Samaka, "SCADA System Testbed for Cybersecurity Research Using Machine Learning Approach," *Future Internet*, vol. 10, no. 8, p. 76, Aug 2018. [Online]. Available: http://dx.doi.org/10.3390/fi10080076

[8] A. H. Mirza, "Computer Network Intrusion Detection Using Various Classifiers and Ensemble Learning," in *2018 26th Signal Processing and Communications Applications Conference (SIU)*, 2018, pp. 1–4.

[9] S. O. Ajakwe, C. I. Nwakanma, D. S. Kim, and J. M. Lee, "Intelligent and Real-Time Smart Card Fraud Detection for Optimized Industrial Decision Process," in *2021 Korean Institute of Communication and Sciences Summer Conference*, vol. 75, 2021, pp. 1368–1370. [Online]. Available: www.dbpia.co.kr/journal/articleDetail?nodeId=NODE10587528

[10] X. Dong, Z. Yu, W. Maharani, Cao, Y. Shi, and Q. Ma, "A Survey on Ensemble Learning," *Frontiers of Computer Science*, vol. 14, no. 8, pp. 241–258, 2020. [Online]. Available: https://doi.org/10.1007/s11704-019-8208-z

[11] C. Yue, L. Wang, D. Wang, R. Duo, and X. Nie, "An Ensemble Intrusion Detection Method for Train Ethernet Consist Network Based on CNN and RNN," *IEEE Access*, vol. 9, pp. 59 527–59 539, 2021.

[12] N. Moustafa, B. Turnbull, and K. R. Choo, "An Ensemble Intrusion Detection Technique Based on Proposed Statistical Flow Features for Protecting Network Traffic of Internet of Things," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 4815–4830, 2019.

[13] A. Abbasi, A. R. Javed, C. Chakraborty, J. Nebhen, W. Zehra, and Z. Jalil, "ElStream: An Ensemble Learning Approach for Concept Drift Detection in Dynamic Social Big Data Stream Learning," *IEEE Access*, vol. 9, pp. 66 408–66 419, 2021.

[14] L. A. C. Ahakonye, C. I. Nwakanma, J. M. Lee, and D. S. Kim, "Evaluating Artificial Intelligence Mitigation Techniques for Countering Attack on Smart Factory SCADA Network," in *2021 2nd Korea Artificial Intelligence Conference (KAI)*, 2021.