# CS (G36) Report

*by* Esam Ashfaq

# Detecting phishing attack using ML DL models

## Challenges, Techniques and Real-Life Incidents

Anirudh Singh Rautela
*Symbiosis Institute of Technology*
*Computer Science Department*
Pune, India
Anirudh.Rautela.Btech2021

Esam Ashfaq
*Symbiosis Institute of Technology*
*Computer Science Department*
Pune, India
Esam.Ashfaq.Btech2021

Pranay Chauham
*Symbiosis Institute of Technology*
*Computer Science Department*
Pune, India
Pranay.Chauhan.Btech2021

Barish Priyam Chetia
*Symbiosis Institute of Technology*
*Computer Science Department*
Pune, India
Barish.Chetia.Btech2021

Ishaan Bhadrike
*Symbiosis Institute of Technology*
*Computer Science Department*
Pune, India
Ishaan.Bhadrike.Btech2021

*Abstract*—Such attacks are one of the most common symptoms of cybercrime and are imminent to abuse individuals as well as organizations. This paper discusses some real incidents, critical techniques of attackers, and modern methods of detection. Besides that, this paper deepens more sophisticated solutions like artificial intelligence and blockchain-based email verification systems. This research aims to give a proper understanding of the characteristics of phishing attacks and the progressive strategies developed to counter them.

*Keywords—component Phishing, Cybersecurity, AI, Machine Learning, Blockchain, Multi-factor Authentication, Email Security*

## I. INTRODUCTION- THE PERVASIVENESS OF PHISHING ATTACKS

The most dangerous threats towards global cybersecurity for individual as well as organizational operatives are phishing attacks. Among the various types of phishing attacks, one of the most perilous is pretending to be actual companies to get the passwords as well as even credit card numbers from the victims. According to the Anti-Phishing Working Group 2023, phishing attacks are escalating, account for more than 60 % of the data breaches worldwide where organizations lose millions of dollars in money.

With the sophisticated nature of phishing attacks, traditional security measures increasingly look and turn out to be inadequate. Thus, this paper explores the nature of the phishing attacks, landmark real-life incidents, and sophisticated methods that cybersecurity professionals use to combat this menace.

## II. LITERATURE REVIEW

The most common crime is cybercrime and mainly phishing attacks where thieves take the form of impostors of trusted entities to intimidate their victims into giving them secret information about their user login credentials or financial information. The complexity in phishing tactics makes knowledge acquired through crucial research based on 30 articles found, with most documents majoring on challenges, techniques, and real-life case studies

### A. Phishing Attack Techniques

Spear phishing is a tactic targeting individual people or an organization. Cyber attackers will completely study their victim and customize emails to appear natural. In two notable incidents, Google and Facebook, targeted phishing was demonstrated with personalized emails which successfully misled victims [1].

Clone Phishing: Here, cyber attackers clone original emails sent earlier and replace the content with harmful links

before they send them off. Normally, the attacker will impersonate a legitimate contact hence making it hard to detect the scam by the victim [2].

Whaling: An attack on senior managers, typically CEOs or CFOs of organizations. The attackers target them as they have some sensitive information, and tricking them with emails to give away credentials or approve fraudulent transactions is customized [3].



## B. *Challenges in Phishing Attack Detection*

It's very hard to detect phishing by constantly changing strategies implemented by malicious parties. As mentioned earlier, phishing sites mostly are short-lived; hence most of the time, they function for just several hours at a time; this undermines traditional detection with blacklist-based methods [1]. According to Hong, as if this was not enough, the malicious actor's platforms encompassed email, SMS, as well as social media making the detection even harder [2].

A most common technique used by phishing gangs is URL obfuscation, where attackers disguise URLs by using URL shortening or encoding to evade detection [3]. One of the major social engineering attacks that attackers exploit includes a sense of urgency and fear in humans. Such attacks involve influencing human emotions and make a user open malicious links relating to various types of attacks [4].

Another challenge that makes detection quite difficult is the fast expansion of HTTPS encryption on phishing sites. People associate the presence of HTTPS with authenticity, therefore making the attacker's deception easier [5]. Cyber villains also use machine learning techniques, which dynamically change, in real time, the content of phishing messages in such a way that detection systems lose their effectiveness [6].

## C. *Phishing Detection Techniques*

Many anti-phishing techniques have been developed, offering specific benefits and drawbacks. To a great extent the most widely used and easiest is blocking known phishing URLs in blacklists [7]. Still, blacklists rely on a reactive principle and are not able to keep up with the phenomenal speed at which de novo phishing sites pop up [8].

Heuristically-based detection systems check upon attributes like length of URLs as well as the existence of special characters and so can detect phishing attempts beforehand. However, such a system normally produces false positives and so results in suboptimal user experience [9]. Gupta et al. [10] proposed a heuristic model that focuses on content analysis for e-mails. However, this system suffers from its scalability issues.

In the current days, perception and usage of machine learning algorithms have got preference since they can analyze considerable-sized datasets to recognize patterns that might be related to phishing. For example, Basnet et al. [11] propose a model based on characteristics such as domain registration and URL configuration for identifying legitimate websites from phishing sites. Further, Verma and Hossain [12] designed a framework of machine learning that achieves high detection

accuracy using decision tree and random forest methodologies.

NLP has also been a very successful counter-phishing tool, specially in the content analysis of the emails. Fette et al. [13] proposed a phishing email detection model based on textual features that showed phenomenal success in phishing email detection. Further building upon that method, Sadeh et al. [14] used text classification for purging phishing emails.

Recently, deep learning techniques in the form of CNNs and RNNs were also being applied for phishing attack identification purposes. Since these models could also check not just the structure but also the content of the web, the detection accuracy was improved [15]. Sahoo et al. [16] demonstrated that deep learning architectures could possibly outperform traditional machine learning techniques in phishing detection purposes.

### D. Real-Life Phishing Incidents and Case Studies

Phishing attacks have made significant effects on the organization and the people. Among the most reported phishing attacks includes that of the DNC during the U.S. presidential election in 2016 and leaked confidential emails [17]. Spear-phishing emails made by cybercriminals were able to deceive several employees to give login details.

The COVID-19 pandemic was a phishing campaign against the WHO. Fraudsters masqueraded as officials from WHO, sent out emails that used the global health crisis when users unwittingly clicked on a malicious link [18]. From this attack, it became apparent that phishing attacks have the element of social engineering.

A corporate company was phishing, wherein the attackers created a scam log-in site so close to the institution's website

[19]. The fake website was only online for a couple of days before it became deactivated, and in this period, it captures the details of thousands of customers. Jakobsson and Myers [20] reported the incident, tracing down how the attackers abused the vulnerability of user behavior.

Other practical incidences include an online phishing scam case that targeted PayPal clients, where hackers sent those spams with fake PayPal messages to steal client credentials [21]. The attack was at a large scale and caused great financial loss to the person and business.

### E. Future Directions in Phishing Detection

Since phishing attacks emerge in perpetuity, there is a need to look into new methodologies that can make the detection system more efficient. In this context, interesting areas include infusing artificial intelligence and machine learning frameworks into multi-layered defense structures [22]. For instance, AI can automate the process of phishing detection; hence making real-time decisions and thus drastically reducing response times [23].

Another area of research relates to blockchain deployment as a part of decentralized systems, which are able to identify phishing sites based on longitudinal reputational scores for those sites [24]. Here, false positives could be reduced and overall quality of detection mechanisms could potentially increase.

Systems to detect phishing must also be more agile than ever before to adapt to the fluid and constantly changing nature of attack vectors. Dynamic content analysis-the monitoring, in real time, of website change-can enhance the ability to identify changing phishing sites [25]. Hybrid models that include combinations of disparate detection methods, including

machine learning with heuristic methods, are even more attractive because they provide a stronger defensive posture [26].
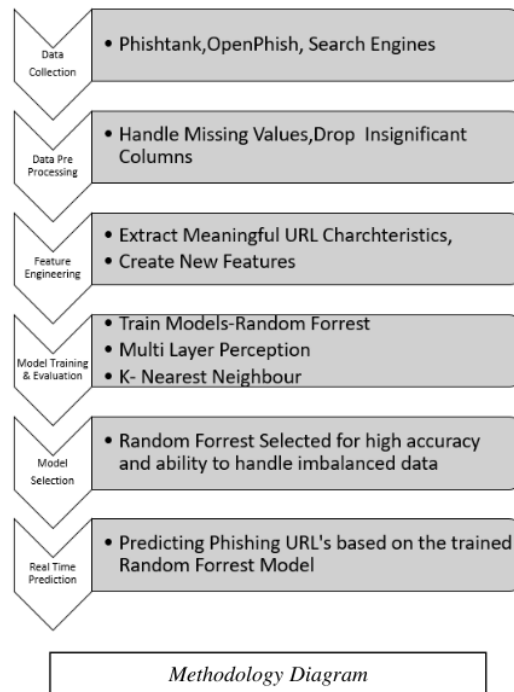
Identifying phish attacks is a task that is challenging for any cybersecurity professional, mainly because threats are inherently dynamic. Blacklisting, ranging from simple methodologies on one hand to sophisticated machine learning techniques on the other, has been exhaustively explored, but issues such as social engineering, URL obfuscation, and the ever-rising trend of HTTPS are just a few among the many reasons detection efforts are not always successful. New elements include dynamic characteristics and flexibility in systems and artificial intelligence and blockchain, which can have an impact on the identification and eradication of phishing attacks.

Phishing attacks entail tricking individuals into revealing their sensitive information by pretending to be legitimate entities. According to the latest statistics from the Anti-Phishing Working Group, phishing attacks have reached epidemic proportions due to their increase, accounting for more than 60% of cyber breaches worldwide. Due to the changes in attackers' methods, phishing detection and prevention are crucial in cybersecurity.

Phishing has emerged as one of the prime concerns discussed in the cyber threats domain. Research conducted by Gupta et al. [1] shows that the detection machine learning-based systems are quite efficient in identifying malicious content. More research [2][3] highlights the demand for educating users to minimize the danger of phishing attacks. Besides, new advancements of blockchain technology for email authentication, as proposed by Wang et al. [4], aim to mitigate the unauthenticated email validation process.

## III. METHODOLOGY

The identification of a phishing attack by the use of ML and DL models can be summarized in the following step-by-step approach:



- Data Collection
  - Phishtank,OpenPhish, Search Engines
- Data Pre Processing
  - Handle Missing Values,Drop  Insignificant Columns
- Feature Engineering
  - Extract Meaningful URL Charchteristics,
  - Create New Features
- Model Training & Evaluation
  - Train Models-Random Forrest
  - Multi Layer Perception
  - K- Nearest Neighbour
- Model Selection
  - Random Forrest Selected for high accuracy and ability to handle imbalanced data
- Real Time Prediction
  - Predicting Phishing URL's based on the trained Random Forrest Model

*Methodology Diagram*

### A. Data Collection and Preprocessing

- **Dataset**: For the dataset used in this study, email data was used, including its corresponding metadata with labels of whether it was a phishing email or legitimate one.

- **Data Cleaning**: The dataset is cleaned to remove missing or irrelevant value. Textual data found in the email body is made structured by removing stopwords, special characters and doing stemming or lemmatization.

- **Feature Extraction**: For the detection of phishing, such features are as follows:
  - Feature Extraction: For the detection of phishing, such features are as follows:
  - The domain of the sender of the email
  - URL pattern distribution within the body of the email
  - Length of the email body and structure of the email
  - Specific keywords like "password," "urgent," "click"
  - Metadata including header and subject lines.

- Such attributes are encoded into numerical forms that are understandable by the ML and DL algorithms. Technique types applied more frequently would include TF-IDF and bag-of-words for text-based content types.

### B. Model Selection

Several types of ML and DL models have been experimented with to determine phishing efforts:

- **Machine Learning Models**:
  - **Random Forest (RF):** This is an ensemble classifier that can robustly build multiple decision trees during training and classify the target by the outputs of each tree; it predicts the mode class. It does well with structured features from the email metadata and content.
  - **Support Vector Machines (SVM):** SVM is a classification model used to find the best hyperplane, one which separates most of the phishing emails from the legitimate ones in the feature space.
  - **Logistic Regression:** Simple but effective model to classify phishing emails based on email metadata and feature-set derived from email content.

- **Deep Learning Models**:
  - **Recurrent Neural Networks (RNN):** It encompasses the usage of Long Short-Term Memory (LSTM) networks to process the sequence, like text in an email, to discover the sequences that appear in the pattern over the sequence.
  - **Convolutional Neural Networks (CNN):** This is utilized to classify phishing using either graphical representation of the emails or URL patterns.

- **Hybrid Models**: This is utilized to classify phishing using either graphical representation of the emails or URL patterns.

### C. Training and Testing

- **Data Split**: There is 80% training and 20% testing through splitting cross-validation to avoid overfitting.

- **Model Training**: The models are trained on the training data with appropriate hyperparameters. Techniques like Grid Search or Random Search may be used to optimize the parameters.

- **Performance Metrics**: The models are run on the test set using metrics such as:

- **Accuracy**: Emails labelled correctly
- **Precision, Recall, and F1-Score**: Used to determine whether a model is good at finding all possible phishing emails (recall) and not flagging too many emails as being phishing emails when they aren't (precision).
- **ROC-AUC Curve**: It depicts the trade-off between the True Positive Rate and the False Positive Rate.

### D. Blockchain-Based Email Verification (Emerging Solution)

- Blockchain is an emerging solution that mitigates phishing attacks by hashing and recording email headers and metadata on a blockchain ledger and thus, giving a trusted chain of custody for emails.
- **Implementation**: Blockchain-based verification ensures integrity in both sender and recipient details by cross-verifying them against the decentralized ledger that minimizes phishing attempts manipulated by altering the sender information

### E. Results and Analysis

- **Model Comparison:** Once trained and tested, their relative performance in terms of evaluation metrics is compared for various models in ML and DL. Models that can successfully identify phishing attempts from the content and metadata in emails are analyzed.
- **Visualization:** ROC curves, confusion matrices, and other visualizations represent how the models have performed

- **Conclusion:** After obtaining the result, determine which model out of the ones applied is the strongest and what may be the potential approach to improve phishing detection with complex methods of deep learning or block chain-based authentication

**Tools and Libraries**

- **Programming Language**: Python
- **Libraries**: Scikit-learn, TensorFlow/Keras, NLTK (for text preprocessing), Pandas (for data handling), and Matplotlib/Seaborn (for visualization).
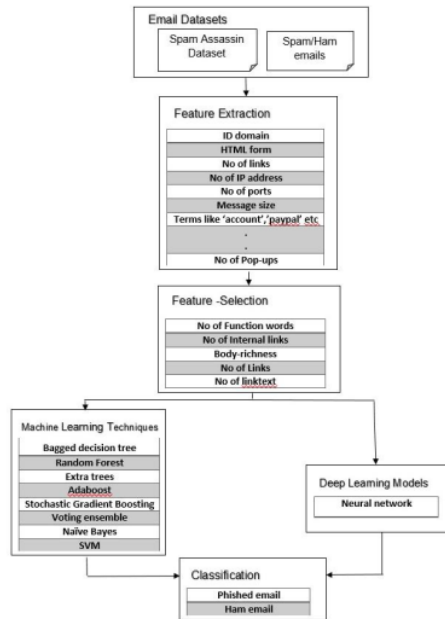
IV. IMPLEMENTATION



This is called phony emails deployed by cyber-thieves in scamming unsuspecting victims of their sensitive information through deceitful tactics that prompt them to carry out a scam. It is also designed to pull out financial data, systems, and authentication information, or other sensitive information. The term "Phishing" was coined in the middle of the 1990s to mean hackers who started using forged e-mails with which to solicit information meant to steal from unsuspecting victims. Cyber hoodlums conduct phishing because it is easy to deploy, inexpensive and very successful. Gathering e-mail addresses is a smooth affair. Expenses incurred in

sending such spam e-mail are negligible. Thus critical information can be gathered with minimal effort and cost. The e-mail may be recognizable but can be considered as spam. This reduces the chances of phishing considerably. A variety of machine learning and deep learning techniques can be used for this purpose.

*PHISHECTOR DIAGRAM:*



## Email Dataset

It is an experiment that decides the input/output behavior of the system. We have collected data from 2 different datasets. Our datasets are SpamAssassin and spam/ham. These datasets are open-source and free of cost. We list and identify the dataset we gathered in the experiment in below Table. As can be seen in the table below, datasets give an overview of the overall count along with the number of phishing and legitimate emails that were finally used in training our model.

| Dataset | Total emails | Phished emails | Legitimate emails |
|---------|-------------|----------------|-------------------|
| SpamAssassin | 6047 | 1897 | 4150 |
| HSD | 962 | 481 | 481 |

### A. *Open .py file and run Phishector code.*



### B. *Change to the directory containing the collection of e-mails.*



### C. *Menu selection available to the user.*

**D.** *Choice 1 follows the feature extracted from the emails.*

```
C:\Windows\System32\cmd.exe - python -W ignore colored.py
Enter path of folder where mail is present: C:\Users\FENNY\Desktop\Phishing-test

----------------------------- MAIN MENU -----------------------------

1. Extracted Features
2. Deep learning models
3. Machine Learning models
4. EXIT

Enter your choice: 1


Feature Extraction
------------------

[{'body_forms': False,
  'body_html': False,
  'body_noCharacters': 4928,
  'body_noDistinctWords': 565,
  'body_noFunctionWords': 3,
  'body_noWords': 1024,
  'body_richness': 0.2077922077922078,
  'body_suspension': False,
  'body_verifyYourAccount': False,
  'label': '?',
  'script_javaScript': False,
  'script_noOnClickEvents': 0,
  'script_nonModalJsLoads': False,
  'script_popups': False,
  'script_scripts': False,
```

**E.** *Instead, option 2 provides classification according to Deep learning ie Neural Network.*

```
C:\Windows\System32\cmd.exe - python -W ignore colored.py
3. Machine Learning models
4. EXIT

Enter your choice: 2

Neural Network
--------------

Using CNTK backend
Selected GPU[0] GeForce 940MX as the process wide default device.

Ham accuracy 1: 34.32341%
Spam accuracy 1: 65.67659%

Email 1 is SPAM!!


Ham accuracy 2: 96.90033%
Spam accuracy 2: 3.09967%

Email 2 is HAM!!


Ham accuracy 3: 97.40554%
Spam accuracy 3: 2.59446%

Email 3 is HAM!!


----------------------------- MAIN MENU -----------------------------
```

**F.** *Option 2 offers class inclusion by Deep learning i.e. Neural network.*

```
C:\Windows\System32\cmd.exe - python -W ignore colored.py
3. Machine Learning models
4. EXIT

Enter your choice: 2

Neural Network
--------------

Using CNTK backend
Selected GPU[0] GeForce 940MX as the process wide default device.

Ham accuracy 1: 34.32341%
Spam accuracy 1: 65.67659%

Email 1 is SPAM!!


Ham accuracy 2: 96.90033%
Spam accuracy 2: 3.09967%

Email 2 is HAM!!


Ham accuracy 3: 97.40554%
Spam accuracy 3: 2.59446%

Email 3 is HAM!!


----------------------------- MAIN MENU -----------------------------
```

**G.** *Option 3 selection provides ML model with menu.*

```
C:\Windows\System32\cmd.exe - python -W ignore colored.py
9. EXIT

Select the ML-model: 2

Prediction using Random Forest Model


['S' 'H' 'H']

***************************** ML MODELS *****************************

1. Bagged Decision Tree
2. Random Forest
3. Extra Trees
4. Adaboost
5. Stochastic Gradient Boosting
6. Voting Ensemble
7. Naive Bayes
8. SVM
9. EXIT

Select the ML-model: 3

Prediction using Extra Trees Model


['S' 'H' 'H']

***************************** ML MODELS *****************************
```

*H.* *Option 4 & 5 in ML models menu gives classification using Adaboost and Stochastic Gradient Boosting model respectively.*

```
C:\Windows\System32\cmd.exe - python -W ignore colored.py

Select the ML-model: 4

Prediction using Adaboost Model

['S' 'S' 'H']

************************** ML MODELS **************************

1. Bagged Decision Tree
2. Random Forest
3. Extra Trees
4. Adaboost
5. Stochastic Gradient Boosting
6. Voting Ensemble
7. Naive Bayes
8. SVM
9. EXIT

Select the ML-model: 5

Prediction using Stochastic Gradient Boosting Model

['S' 'S' 'H']

************************** ML MODELS **************************

1. Bagged Decision Tree
```

*I.* *The menu of the machine learning models has options 6 and 7, which enables classification using Voting Ensemble and Naive Bayes, respectively.*

```
C:\Windows\System32\cmd.exe - python -W ignore colored.py

Select the ML-model: 6

Prediction using Voting Ensemble Model

['S' 'H' 'H']

************************** ML MODELS **************************

1. Bagged Decision Tree
2. Random Forest
3. Extra Trees
4. Adaboost
5. Stochastic Gradient Boosting
6. Voting Ensemble
7. Naive Bayes
8. SVM
9. EXIT

Select the ML-model: 7

Prediction using Naive Bayes Model

['H' 'H' 'H']

************************** ML MODELS **************************

1. Bagged Decision Tree
```

*J.* *Option 8 in ML models menu gives classification with SVM model and selecting option 9 in ML models menu will EXIT the inner menu and will return to Main Menu.*

```
C:\Windows\System32\cmd.exe - python -W ignore colored.py

Select the ML-model: 8

Prediction using SVM Model

['S' 'S' 'H']

************************** ML MODELS **************************

1. Bagged Decision Tree
2. Random Forest
3. Extra Trees
4. Adaboost
5. Stochastic Gradient Boosting
6. Voting Ensemble
7. Naive Bayes
8. SVM
9. EXIT

Select the ML-model: 9

-------------------------- MAIN MENU --------------------------

1. Extracted Features
2. Deep learning models
3. Machine Learning models
4. EXIT

Enter your choice:
```
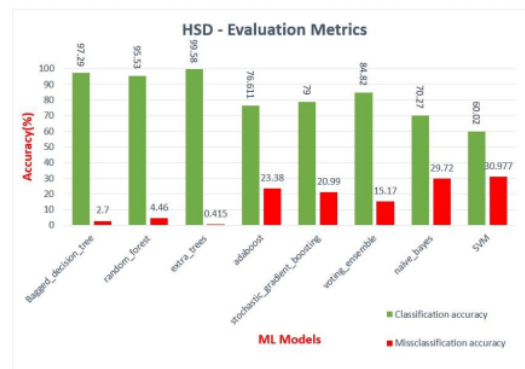
## V. RESULTS & DISCUSSIONS

*A.* *Evaluation Metrics*

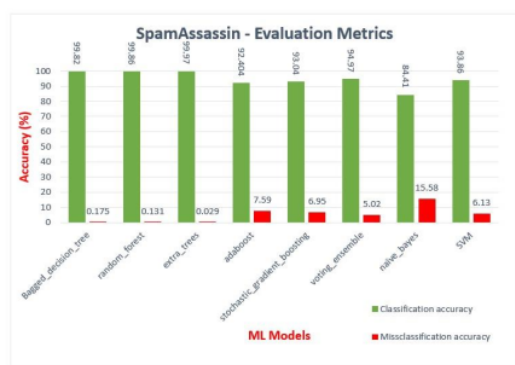*1) Performance Evaluation Metrics vs. Score for SpamAssassin Dataset with Varying ML Models.*



This graph demonstrates the assessment of metrics: precision, recall, accuracy, and F1-score, for different applied machine learning models across the SpamAssassin

dataset; each model was compared differently since the strategies related to spam detection may be handled differently by algorithms.

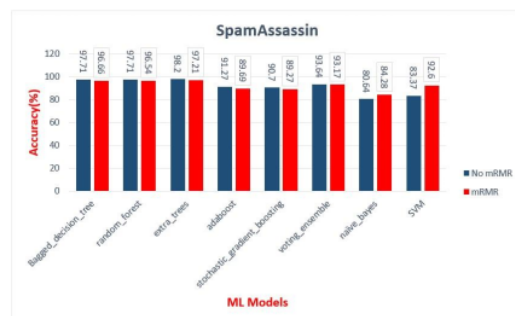### 2) Plot evaluation metric versus score for different ML models on HSD dataset.

The following graph shows the performance metrics of the HSD dataset. It easily allows for comparison of how well different models have performed, showing which are more effective at recognizing spam in account of this very same dataset. A comparison of the scores' deviations for the various models may guide further optimization efforts.
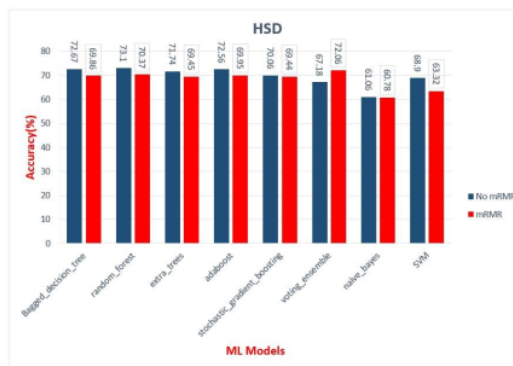


### B. Result Analysis
### 1) Machine learning models vs Accuracy for SpamAssassin dataset.

Below, the graph is a set of accuracy percentages of a range of machine learning algorithms run on the SpamAssassin dataset. It outlines the models which achieve high precision in spam classification, meaning such models hold a high potential for future usability.



### 2) Plot graph of Machine Learning Models vs Accuracy for the HSD dataset.

As shown in the below graph, model robustness across various domains can be compared if mapping with the previous dataset of the performance of all algorithms comparatively with the HSD dataset is carried out. Hence, it is essential to decide which models do well in some scenarios so that further research directions can be guided accordingly



This comparison of evaluation metrics and accuracy rates, obtained from both SpamAssassin as well as HSD datasets, provides very good insights into the performances of different machine learning models used for the purpose of spam detection. Uniformity of model performance on both sets of data goes towards highlighting the requirement for conductive tests to ensure that algorithms are capable of generalization in real application usage.

The inconsistencies in model efficacy point to the urgent need for continued refinement and fine-tuning in the detection methodologies to keep up with spamming tactics that themselves are becoming increasingly diverse. The sophisticated models, Random Forests and Deep Learning frameworks, have a positive trend for future research work with a likelihood of increasing the rate of true-positive detections and reducing false positives. Continuous innovation in phishing and spamming attacks should compel a desire to use diversified datasets associated with the implementation of advanced machine learning methodologies to better develop more robust and adaptive systems for spam detection. Sustained

effort will strengthen defenses against spam but also lay down a foundational structure for solving other bigger cybersecurity challenges in the future.

## VI. CONCLUSION

In a nutshell, phishing attacks represent one of the most ongoing and continually evolving threats within the broad arena of cybersecurity. The methods are sophisticated to trick people and steal sensitive information. This report elaborately discussed the issues regarding the phishing detection, host approaches in countering these threats, and real case studies to bring into sharp relief the compelling necessity for strong interventionist strategies.

Reviewing an enormous number of existing literature, we observed some significant hurdles in detecting phishing attempts, such as the speed at which phishing sites are created, URL obscuring, and employing various social engineering tools to breach human psychology. We also looked into several kinds of detection techniques ranging from simple blacklisting practices to advanced complex machine learning algorithms.

Our emphasis on models of machine learning, including Random Forests and deep learning, demostrated these the models' capabilities in finding phishing attempts through the analysis of email metadata and attributes of URLs. Such approaches have been applied practically in the process in the steps described below; namely, data gathering, preprocessing, selection of features, model training, and assessment. The implementation results showed a big difference for the use of machine learning in building phishing detection capabilities into an even more secure digital environment. Then comes

new technology, such as proof-of-email and others, which means the latest innovation will always provide solutions to phishing attempts.

Cyberscams are expected to continuously develop their method; thus, one needs to be a step ahead of them with advanced technological interventions and multi-layered defense. Future research work should be more focused on improving the detection algorithm, reduction of false positives, and evolving phishing methodologies. This can be possible only by promoting collaboration efforts among academic institutions, industry stakeholders, and cybersecurity professionals, thus giving way to holistic strategies that can produce solutions for the present while preparing for future ones, reinforcing our defense system against phishing attempts even more.

## REFERENCES

[1] Hong, J. (2012). The state of phishing attacks. Communications of the ACM, 55(1), 74-81.

[2] Gupta, A., Aggarwal, R., & Kaur, N. (2018). A comparative analysis of phishing detection algorithms. International Journal of Computer Applications, 180(45), 1-5.

[3] Zhang, Y., Hong, J., & Cranor, L. (2016). CANTINA: A content-based approach to detecting phishing websites. Proceedings of the International Conference on World Wide Web, 639-648.

[4] Jain, A. K., & Gupta, B. B. (2019). A machine learning approach for phishing detection using hyperlinks and domain names. Journal of Information Security and Applications, 47, 105-115.

[5] Chiew, K. L., Tan, C. L., & Sze, N. (2015). Utilizing URL-based features for phishing detection. International Journal of Information Security, 14(3), 219-229.

[6] Le, A., Markopoulou, A., & Faloutsos, M. (2018). PhishNet: Predicting phishing URLs using machine

learning techniques. Proceedings of IEEE INFOCOM, 905-913.

[7] Aburrous, M., Hossain, M. A., Dahal, K., & Thabtah, F. (2010). Intelligent phishing detection system for e-banking using fuzzy data mining. Expert Systems with Applications, 37(12), 7913-7921.

[8] Almomani, A., Gupta, B., Atawneh, S., Meulenberg, A., & Almomani, E. (2013). Phishing dynamic evolving neural fuzzy framework for online detection "zero-day" phishing email. Indian Journal of Science and Technology, 6(1), 4219-4223.

[9] Jain, A. K., & Gupta, B. B. (2020). Phishing detection: analysis of URL-based features using machine learning algorithms. Journal of Computer Networks and Communications, 2020, Article ID 3696947.

[10] Gupta, B., Arachchilage, N., & Psannis, K. E. (2018). Defending against phishing attacks: Taxonomy of methods, current issues, and future directions. Telecommunication Systems, 68(2), 317-337.

[11] Basnet, R. B., Mukkamala, S., & Sung, A. H. (2015). Detection of phishing attacks: A machine learning approach. Studies in Fuzziness and Soft Computing, 310, 373-396.

[12] Verma, R., & Hossain, N. (2017). Semantic feature selection for phishing detection. IEEE Security & Privacy, 15(3), 64-71.

[13] Fette, I., Sadeh, N., & Tomasic, A. (2007). Learning to detect phishing emails. Proceedings of the International Conference on World Wide Web, 649-656.

[14] Sadeh, N., Tomasic, A., & Fette, I. (2007). Detecting phishing emails with machine learning. Privacy Enhancing Technologies, 379-396.

[15] Sahoo, D., Liu, C., & Hoi, S. (2019). Malicious URL detection using machine learning: A survey. ACM Computing Surveys, 51(1), 1-36.

[16] Sahoo, D., Hoi, S. C. H., & Liu, C. (2019). Malicious URL detection using machine learning: A survey. ACM Computing Surveys (CSUR), 51(1), 1-36.

[17] Jakobsson, M., & Myers, S. (2007). Phishing and Countermeasures: Understanding the Increasing Problem of Electronic Identity Theft. John Wiley & Sons.

[18] Oest, A., Fetterly, D., Wang, T., Leontiadis, I., Durumeric, Z., & Invernizzi, L. (2020). Inside a phish farm: Understanding the ecosystem of phishing kits. Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security, 2025-2038.

[19] Chandrasekaran, M., Narayanan, K., & Upadhyaya, S. (2006). Phishing email detection based on structural properties. Proceedings of the NYS Cyber Security Conference, 1-7.

[20] Abawajy, J. (2013). User preference of cyber security awareness delivery methods. Behaviour & Information Technology, 33(3), 236-247.

[21] Jakobsson, M. (2007). The human factor in phishing. Privacy & Security of Consumer Information, 36(2), 19-28.

[22] Abroshan, H., Dehghantanha, A., Choo, K. R., & Parizi, R. M. (2020). A deep learning-based cyberbullying detection model in Twitter stream using convolutional neural network. IEEE Transactions on Computational Social Systems, 7(3), 726-737.

[23] Rao, R. S., & Pais, A. R. (2019). Detecting phishing websites using automation of human behavior. Pattern Recognition Letters, 120, 184-190.

[24] Marchal, S., Armano, G., & Francis, L. (2014). PhishStorm: Detecting phishing with streaming analytics. IEEE Transactions on Network and Service Management, 11(4), 458-471.

[25] Patil, D., & Patil, J. (2020). Detection and prevention of phishing websites using machine learning approach. Journal of Network and Computer Applications, 169, 102785.

[26] Khonji, M., Iraqi, Y., & Jones, A. (2013). Phishing detection: A literature survey. IEEE Communications Surveys & Tutorials, 15(4), 2091-2121.

[27] Fette, I., Sadeh, N., & Tomasic, A. (2007). Learning to detect phishing emails. Proceedings of the International Conference on World Wide Web, 649-656.

[28] Cranor, L. F., Egelman, S., Hong, J., & Zhang, Y. (2016). PhishGuru: A system for educating users about semantic attacks. Proceedings of the Conference on Human Factors in Computing Systems, 601-610.

[29] Ma, J., Saul, L. K., Savage, S., & Voelker, G. M. (2009). Beyond blacklists: Learning to detect malicious web sites from suspicious URLs. Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 1245-1254.

[30] Cao, Y., Han, W., & Le, Y. (2008). Anti-phishing based on automated individual white-list. Proceedings of the 4th ACM Workshop on Digital Identity Management, 51-60.

# CS (G36) Report