# Leveraging Large Language Models (LLM) for the Detection of Spear-Phishing Emails as Indicators of Advanced Persistent Threats (APTs)

Aslam Malik Abdul Azeez

Student ID: x23183098

School of Computing

National College of Ireland

Supervisor: Mr. Imran Khan

| **Student Name:** | Aslam Malik Abdul Azeez | | |
|---|---|---|---|
| **Student ID:** | X23183098 | | |
| **Programme:** | MSc in Cybersecurity | **Year:** | 2024-2025 |
| **Module:** | Practicum | | |
| **Supervisor:** | Mr. Imran Khan | | |
| **Submission Due Date:** | 29/01/2025 | | |
| **Project Title:** | **Leveraging Large Language Models (LLM) for the Detection of Spear-phishing Emails as Indicators of Advanced Persistent Threats (APTs)** | | |

**Word Count:** 8485    **Page Count: 22**

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.
ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| **Signature:** | Aslam Malik Abdul Azeez |
|---|---|
| **Date:** | 29/01/2025 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies) | ☐ |
| **Attach a Moodle submission receipt of the online project submission,** to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Leveraging Large Language Models (LLM) for the Detection of Spear-Phishing Emails as Indicators of Advanced Persistent Threats (APTs)

Aslam Malik Abdul Azeez

X23183098

**Abstract**

Spear phishing and Advanced Persistent Threats (APTs) are targeted and context specific, they escape detection by traditional systems (Xuan, 2021). In this research, an advanced detection framework is developed using state-of-the-art machine learning (ML) techniques. The system extracts feature from the content of email and (Innab et al., 2024) email header and social behaviour data to identify language anomalies, metadata patterns and user activity profiles indicative of threats.

The framework reaches high accuracy, precision, recall, F1 scores using ML models such as deep learning and supervised learning, surpassing traditional systems. More specifically this work focuses on the utilization of ML methods to mitigate cybersecurity risks and adds to the burgeoning field of intelligent threat detection systems (Innab et al., 2024) which pinpoints the importance of data in improving organizational security. However even more work needs to be done, but this approach is a promising step forward to counter cyber threats.

## 1     Introduction

Spear phishing and Advanced Persistent Threats (APTs) have matured into cyber threats of increasing sophistication that circumvent traditional detection methods with purpose built cybergangs targeting specific individuals or organizations. Spear phishing uses a deceptive email but an APT is a more stealthy, longer lasting activity to compromise systems or obtain stolen data as they come. Often these adaptive and stealthy threats defeat traditional security solutions, which rely on predefined rules. Machine learning (ML), but especially supervised learning, can discover powerful, malicious patterns in massive amounts of data instead of with static rules. ML models can be trained to look at email metadata, content and sender behaviour for deception for spear phishing. An analysis of user behaviour, network activity is used to identify anomalies and potentially malicious activity for Pts. This relies on data complex systems with no guarantee that they will identify and respond to evolving cyber threats, but they offer a promising, adaptive means of doing so.

## 2     Research Question

" How can machine learning techniques be leveraged to improve the detection of spear phishing and Advanced Persistent Threats (APT) using email analysis and behavioural patterns?"

# 3    Research Motivation & Objective

Conventional security systems are struggling to compete with the swift progression of cyber threats, spear phishing and Advanced Persistent Threats (APT). For many of these attacks they are extremely targeted, stealthy and developed to exploit specific individuals or organizations, and these types of attacks are hard to catch with traditional knowledge. Spear phishing is about misleading email communications, while APTs are about elaborate, long- drawn-out attacks designed to steal data. However, both point to some of the most pressing cybersecurity threats in the present. These new and evolving threats clearly make traditional signature-based systems increasingly ineffective in efforts to prevent attacks that behave differently than they were expected to. The need to produce more adaptive, precise and efficient detection mechanisms that can properly defend against these sophisticated attacks (Xuan, 2021) (Innab et al., 2024) drives this research.

Both spear phishing and advanced persistent threats (APTs) are continuing to evolve, which means the strategies used to detect them must as well. Machine learning (ML) is a good idea, promising systems the ability to look through historical data, discover latent patterns, and predict heretofore unknown attacks. This research attempts to support cybersecurity measures using ML techniques, such as supervised learning, to increase detection accuracy, as well as to reduce the odds of successful cyberattacks (Innab et al., 2024) (Shaukat et al., 2020).

One additional factor that must be considered is that organizations have been depending on emails and digital interactions more and more. One reason spear phishing often takes e-mail form is that it is possible to look for indicators of a spear phishing attempt in various characteristics of an e-mail including metadata, or to information and message content. In addition, behavioural analysis (looking at user interaction), is also useful for identifying advanced persistent threats (APTs) by deducing when there is deviation from normal behaviour (Basit et al., 2020).

## Research Objectives

This work aims to explore and implement machine learning techniques aimed at detecting spear phishing as well as advanced persistent threat (APT) attacks, taking email analysis and behavioural patterns as the main data sources. Analysis of email content and metadata, and user and sender behaviours, as well as the interaction of users with emails, to find signals of malicious activity are done (Zhang & Wang, 2023). Additionally, we need a feature engineering method that transforms raw email data and user behaviour to pertinent and significant features for machine learning model training. However, for improving the accuracy of the models, good feature engineering is required (Ding et al., 2021). Furthermore, we should evaluate and compare different machine learning models (e.g. decision trees, support vector machines and ensemble methods) to find the most suitable techniques of detecting spear phishing and Pts. Key performance metrics like accuracy, precision, recall and F1-score (Chawla et al., 2002) will be evaluated as will be the evaluation.

# 4    Literature Review

The identification of spear phishing and Advanced Persistent Threats (APT) has emerged as a crucial focal point in the realm of cybersecurity research. These types of attacks present significant challenges to detection, primarily due to their adaptive characteristics and their propensity to target specific individuals or organizations. Traditional methods of detection— such as those based on signatures or heuristics—often prove inadequate in recognizing these evolving threats. Because of this, machine learning (ML) techniques have garnered substantial interest, given their ability to discern patterns within extensive datasets and to detect previously unidentified attacks.

## Spear Phishing Detection

Spear phishing is a special, targeted type of phishing in which email content, including personalized email messages, is specially tailored to trick the targets. There have been numerous studies that employ machine learning based techniques to detect phishing emails among which most have contributed to both content analysis and metadata characteristics. For example, Inna et al. (2024) used an ensemble learning to spear phishing by embedding content-based features. In their methodology they combined different models to improve the detection accuracy in analysing different aspects of email content including linguistic, header and structural properties of the email (Xuan, 2021). This approach demonstrated improved accuracy, relies mostly on email content and struggled to adapt when attacker uses obfuscation techniques.

In Basit et al. (2021), email metadata from sender's domain and subject line are explored to understand phishing email detection. The research shown here points out that it is important to extract several features (e.g., time patterns, IP addresses, domain reputation) to effectively identify malicious email. Although based on this method, the efficacy in recognizing basic phishing attempts was demonstrated, it failed to correctly identify more sophisticated attacks with the use of social engineering tactics (Innab et al., 2024). This approach limits its applicability to advanced scenarios.

According to Shaukat and Luo (2020), a machine learning based system for detecting phishing emails is proposed by them using a hybrid approach, where content and metadata analyses are combined. Moreover, for classification, decision trees and support vector machines (SVM) were used in this methodology. The integration of these features increased accuracy according to their findings but disregards behavioural analysis an indispensable component in distinguishing more advanced persistent threat (APT) attacks. Static features still relied on the study and were not used with a continuous learning mechanism to adapt to constantly evolving phishing tactics (Basit et al., 2020).

## APT Detection

Advanced Persistent Threats are defined by their sustained and covert nature. To truly identify APTs, they should be analysed across multiple sources of data including network activity, user behaviours and system logs. Current methods of detecting APTs are based on signature or heuristic methods but may return to poor APT performance once new APT tactics are

employed. To discover APTs, Xuan et al. (2021) used machine learning to do anomaly detection. Instead, they integrated network traffic data with user behaviour to underline deviations from normal behaviour. For identifying potential indicators of APTs, the researchers carried out clustering algorithms and outlier detection techniques. They succeeded in showing effectiveness in uncovering unknown APTs, using predefined baselines of what constitutes 'normal' behaviour; however, a high rate of false positives was shown (Innab et al., 2024) as updating these baselines frequently becomes difficult.

Zhang et al. (2023) proposed a novel method for feature engineering of machine learning and time series analysis of network data to discover Advanced Persistent Threats (APTs). They showed that analysing long term data and user behaviour patterns has helped in improving APT detection. Nevertheless, they had limitations in terms of management of large data volumes, while meeting an immediate requirement for real time detection as the model efficiency decreased with complexity of the network traffic (Shaukat et al., 2020). Zhang et al. (2023) in a related study enact a hybrid model which combining multiple sources of data including network traffic and email metadata to more effectively detect APTs. Due to the high accuracy needed for detection, this methodology employed both supervised and unsupervised learning techniques to improve detection accuracy by six.

## Machine Learning Approaches and Methodologies

Phishing and APT (Advanced Persistent Threat) detection systems that employ Machine learning techniques have shown very promising potential. It possesses an independent ability to learn patterns from data and thereby becomes more adaptable and more adept at recognizing previously unidentified threats. In this domain, some popular ML algorithms used are decision trees, support vector machines (SVM) and ensemble learning methods. It is because decision trees and SVMs are often used for phishing detection because they are easy to interpret and perform well with smaller datasets. Shaukat and Luo (2020) proposed the combination of decision trees and SVMs for phishing detection, with some significant accuracy, but also warned of issues facing the detection against dynamic and real-time threats. While both methods are sufficient for basic classification tasks, they may fail to deal with complexity of APTs (Basit et al., 2020) where feature sets are not known as a priori. Its reliance can however hinder their performance when the landscapes are evolving fast.

Ensemble learning (combining the outputs of many models) requires little modification to obtain significant gains in spear phishing attack detection, according to Inna et al. (2024), as it synergistically combines multiple classifiers using different criteria for analysing the data. Although this method is quite a robust approach, it also might have an impact on computational costs and speeds of detection (Xuan, 2021). Advanced persistent threats (APTs)have commonly been clustered and anomaly detection techniques to identify. Xuan et al. (2021) used clustering algorithms to derive APTs from network activity and user behaviour. However, these techniques depend highly on what constitutes normal and can then predict deviations from it. That dependence can mean false positives in quickly changing environments (Innab etal., 2024).

**Efficiency and Improvements in Current Work**

There are two major impediments to research in the current day: detecting in real time and adapting to changing attack methods. Many of the existing models make strong use of pre-defined features (or static baselines) that may not be robust against novel attack vectors. Yet however machine learning models surpass traditional methods, they have issues like higher ratios of false positives and huge creation expectation of processing immaterial datasets in time. To address these challenges, this proposed work would like to take email analysis and behavioural patterns and use them to find spear phishing and advanced persistent threats (APTs).

# 5     Salient Features

**Multi-Source Data Integration**

There are two major impediments to research in the current day: detecting in real time and adapting to changing attack methods. Many of the existing models make strong use of pre-defined features (or static baselines) that may not be robust against novel attack vectors. Yet however machine learning models surpass traditional methods, they have issues like higher ratios of false positives and huge creation expectation of processing immaterial datasets in time. To address these challenges, this proposed work would like to take email analysis and behavioural patterns and use them to find spear phishing and advanced persistent threats (APTs).

**Behavioural Pattern Analysis and Machine Learning Techniques**

Unlike previous studies that have limited study within one dimension (i.e. email content space or network activity), it is expected that this combination of multiple data sources (such as email metadata, email content analysis, and user behaviour) will yield a more comprehensive and flexible detection system. Additionally, the model's continual learning and updating will help with coping with the problem of changing threats.

The feature vectors are used to categorize attacks generated by emails using their content, metadata and user behaviour using SVMs (Support Vector Machines). In situations where benign and malicious cases separate clearly, they are great. One of their advantages is that they work with high dimensional data (this is important) and they can find patterns for high dimensional feature spaces. We interpret and classify phishing and APT activities using easily comprehensible feature sets using decision trees. These models have rules behind classifications that can be visualized the decision-making process that comes with them.

BERT (also known as Bidirectional Encoder Representations from Transformers) is highly useful when you're trying to understand the finer points and connections within words in emails. Unlike other bag-of-words models, BERT is trained to understand the context in which a word occurs in a sentence, making it excellent at detecting subtle spear phishing (which are often very verbally encoded in natural language). But this type of approach greatly improves

the model's capability to catch complex and contextually adept phishing attempts that can be quite hard to discern using more rudimentary rule based or commonly used machine learning methods.

Like how machine learning algorithms work in picking patterns out of enormous and tangled information, neural systems work, as well (Xuan, 2021). In fact, these networks (namely, deep learning models) can autonomously learn high level features from raw data and the cross section of these features proved highly beneficial for identifying subtle patterns associated with phishing or advanced persistent threat (APT) types of attacks (Innab et al., 2024) (Innab et al., 2024) and detecting new and evolving attack strategies. At the same time, interoperability and deployment of these models can be problematic due to the complexity of these models. (Zhang & Wang, 2023) (Ding et al., 2021).

## Feature Engineering and Data Preprocessing

The success of any machine learning model lies in feature engineering. This research carefully designed the feature extraction process which identifies indicators of phishing and APTs from different raw data sources. To further filter these hypotheses and identify suspicious patterns suggestive of a phishing attempt, we extract email metadata, including the sender IP address (domain) and subject line. However, even email content is important, which are also analysed through natural language processing (NLP) techniques that scan emails content, looking for specific patterns, keywords and linguistic clues to indicate phishing attempts.

Because BERT really works (in identifying the contextual meaning of words in the email) so well, it excels at detecting more subtle and sophisticated phishing tactics. Features of user behaviour (user interaction such as login time and frequency of access of specific systems) are used to detect anomalous patterns associated with possible ongoing advanced persistent threats (APT).The model incorporated network traffic data which includes packet analysis and communication pattern information and through analysing the network behaviour and looking for deviations from normal communication pattern, the model is capable to detect the potential APT might be exfiltrating data or to establish persistent connection. (Shaukat et al., 2020) extracted to help identify suspicious patterns indicative of phishing attempts. Email content, however, is also critical; natural language processing (NLP) techniques are employed to analyse the emails' content, searching for specific patterns, keywords, linguistic cues that are typical of phishing attempts.

Network traffic data, including packet analysis and communication patterns, incorporated into the model; by analysing network behaviour and looking for deviations from typical communication patterns, the model can detect potential APTs attempting to exfiltrate data or establish persistent connection (Shaukat et al., 2020).

Email metadata—such as the sender's IP address (domain) and subject line—are extracted to help identify suspicious patterns indicative of phishing attempts. Email content, however, is also critical; natural language processing (NLP) techniques are employed to analyse the emails' content, searching for specific patterns, keywords, linguistic cues that are typical of phishing attempts.

BERT (Bidirectional Encoder Representations from Transformers) is particularly effective (because it captures the contextual meaning of words within the email), making it adept at identifying more subtle and sophisticated phishing strategies. User behaviour features—related to user interactions (like login time and frequency of access to certain systems)—are used to detect anomalous patterns that may indicate an ongoing advanced persistent threat (APT).

Network traffic data, including packet analysis and communication patterns, incorporated into the model; by analysing network behaviour and looking for deviations from typical communication patterns, the model can detect potential APTs attempting to exfiltrate data or establish persistent connection (Shaukat et al., 2020).

- Phishing Email Datasets: These datasets (which contain many phishing and legitimate emails) allow the model to learn how to discern between malicious and benign content. They encompass features like email body text, sender information and metadata—essential elements for phishing detection tasks.
- APT Datasets: Datasets related to APT (Advanced Persistent Threat) attacks often consist of system logs, user behaviour data and network traffic logs. This data gives us insights on how attackers act and over time help us model anomalous patterns that can stem from long term persistent threats.
- Custom Behavioural Datasets: In addition, custom datasets were also captured to reflect the real time user interaction patterns and network traffic captured with the publicly available datasets.

## Evaluation and Performance Metrics

The model's performance is rigorously evaluated using standard machine learning metrics (e.g., accuracy, precision, recall, F1-score). These metrics are essential for assessing effectiveness (of the detection system) in identifying both phishing emails and APT activities. The evaluation process—however—includes cross-validation and testing on unseen data, because it ensures model robustness and ability to generalize to new threats (Zhang & Wang, 2023).

# 6    Dataset and Machine Learning Mechanism

## Datasets Used in Research

Machine learning model efficacy is strongly correlated to the data quality and diversity used for training and evaluation (Xuan, 2021). The datasets play a vital role in the performance of the detection system, in this research. Employed for training and evaluation. In this research, datasets assume a crucial function is the effectiveness of the detection system. These primary datasets include phishing email datasets, APT related datasets and customized behavioural datasets developed to detect unusual user behaviour. APT-related datasets and customized behavioural datasets designed recognize atypical user activity. Each dataset is chosen to ensure the identification of both spear phishing emails and APT-related behaviours, thus establishing robust foundation for training the machine learning models. But the selection process is a complex one, since it demands,

due attention be paid to many factors (Innab et al.,2024).

Phishing email datasets are labelled instances of phishing emails and legitimate emails. The training for these models to identify spear phishing attempts is based on email content and metadata; these datasets serve that purpose. Phishing Email Dataset (KDD Cup), Enron Email Dataset, Spam Assassin Dataset are phishing email datasets publicly available examples. This research: the datasets assume a crucial function the effectiveness of the detection system. The primary datasets utilized include phishing email datasets, APT-related datasets and customized behavioural datasets designed recognize atypical user activity. Each dataset is chosen to ensure the identification of both spear phishing emails and APT-related behaviours, thus establishing robust foundation for training the machine learning models. However, the selection process can be complex, because it requires careful consideration of various factors (Innab et al., 2024).

Phishing email datasets encompass labelled instances of both phishing and legitimate emails. These datasets serve the purpose of training models aimed at identifying spear phishing attempts, which rely on email content and metadata. Publicly accessible examples phishing email datasets are the Phishing Email Dataset (KDD Cup), Enron Email Dataset and Spam Assassin Dataset. However, the selection process can be complex, because it requires careful consideration of various factors (Innab et al., 2024). Phishing email datasets encompass labelled instances of both phishing and legitimate emails. These datasets serve the purpose of training models aimed at identifying spear phishing attempts, which rely on email content and metadata. Publicly accessible examples phishing email datasets are the Phishing Email Dataset (KDD Cup), Enron Email Dataset and Spam Assassin Dataset. Typically, these datasets incorporate features such as email subject lines, sender details and message content. The Phishing Email Dataset is especially beneficial; this is primarily because it comprises labelled emails that span a diverse array of attack types. The machine learning model can therefore assimilate patterns embedded in phishing and legitimate emails (Xuan, 2021).

## Machine Learning Models Used

The Support Vector Machines (SVM) are extensively used classification algorithms in the realm of Machine Learning that are known for being able to find out the best decision boundaries in the high dimensional space. The categorization of phishing and APT (Advanced Persistent Threat) attacks is performed using SVMs on feature vectors constructed from email content, metadata and user behaviour (Xuan, 2021). Given the clear margin of separation, SVMs perform very well distinguishing (and separating) between 'malicious' and 'benign' instances.

In this research, decision trees are used to represent important models in machine learning. Clear, interpretable, and a method to comprehend classification decisions, but building a tree like structure, each internal node represents a feature, each leaf node represents a decision. With their good handling of categorical and numerical data, decision trees are especially good at classifying phishing emails and identifying APTs (Innab et al., 2024) (Innab et al., 2024).

In this research we focus on such decision trees as models. They provide clear and interpretable methods for understanding classification decisions: a tree in which each internal

node is a feature, and each leaf node is a decision, forms a tree like structure. Decision trees are great    choices when classifying, phishing emails and identifying APTs because they can handle categorical as well as numerical data. In addition to that, they simplify the reasoning behind each of the decisions, giving a valuable overview of the attack detection process (Zhang & Wang, 2023).

The model that analyses part of the content of mails is known as BERT. BERT differs from traditional word-based approaches that rely on treating words in isolation, as opposed to colours pervasive knowledge of words from their surrounding text via contextual embeddings. The ability to notice more subtle details such as intent behind a phishing email, something regular email client cannot do. While BERT is very good at spotting complex, context dependent spear phishing techniques that aren't just based on word hybrids, it's also an important tool to find out about advanced phishing schemes. Subsequently, the model was trained with labels of simulated phishing and normal emails, to define fine grain patterns indicative of malicious intent (Chawla et al., 2002).

The autonomous learning of features from raw data without feature engineering by neural networks (especially deep learning models) is utilized in this research to identify evasive attack patterns that are not apparent to the eye. We apply the deep learning methodology to both email content associated with user behaviour data, which lets the system detect anomalies in very large datasets (Chawla et al., 2002).

## Detection Mechanisms

The system starts off by scrutinizing email content and metadata for the common signs of phishing attempts (initially). It is essentially looking at suspicious senders, unusual subject lines, unusual patterns in email content. The email is interpreted utilizing Natural Language Processing techniques (BERT) to understand the semantic meaning of the email to aid detecting more sophisticated phishing strategies based on contextual manipulation. In addition, the system monitors user behaviour to detect anomalies indicative of an Advanced Persistent Threat (APT) such as user login times, data access and network resource interaction habits. If sometimes the user accesses uncommon files or system, it will mark down as attack used. Common in APT attacks, the analysis of network traffic (notably) also looks for data exfiltration or lateral movement indications. Since the pattern is rarely the same, the system learns to notice individual setups, which might be related to malicious activities.

## Feature Engineering and Data Preprocessing

Indeed, feature engineering helps with effective feature engineering (for model performance). In this work we employ various methods to derive meaningful features from raw data in the form of email text, user interaction logs and network traffic. It features the important phrases statistics (linguistic indicators) and tokenized words out of the email content extraction methods, for example. BERT embeddings are useful for capturing contextual relationship between words and phrases but to create user behaviour profiles, other facts have been collected from system logs such as login times, time spent on resources accessed by the user and user action.  Network features: Based on network traffic logs, elements such as packet sizes, destination IP addresses and protocol types are extracted to detect suspicious activities. Cleaning and transforming raw

data, applying dimensionality reduction to reduce model complexity or numerical features, while normalizing numerical features for easier comparison are all preprocessing steps (Shaukat et al., 2020).

# 7      Research Methodology

## Problem Identification and Definition

The research starts by pinpointing the specific issue at hand: identifying spear phishing and APT attacks within a networked setting. These types of attacks tend to be subtle and complex, utilizing methods like social engineering for spear phishing and extended, unnoticed infiltration for Pts. The challenge is to develop a system capable of automatically detecting suspicious emails and unusual network or user activities that could signal these attacks.

To solve this problem, the research focuses on two main objectives:

- **Phishing Email Detection**: Automatically identifying emails that contain phishing attempts, particularly spear phishing, which is targeted at specific individuals or organizations using social engineering techniques.
- **APT Detection**: Identifying malicious activities associated with APTs, such as abnormal network traffic, unauthorized system access, or unusual user behaviour.

## Detection Mechanisms and Models

The detection mechanism is a multi-layered approach that combines different models and algorithms to identify phishing and APT attacks from various data sources. The primary components of the methodology include:

## Email Content Analysis (Phishing Detection)

• They gather a very large dataset of phishing and legitimate emails sourced from the Phishing Email Dataset and the Enron Email Dataset and custom datasets from the target environment.

• Cleaned emails removes unnecessary data (egg stop words and special characters) and processed with tokenization, stemming, or lemmatization to reduce the text.

• These features include keyword frequency, sender reputation, subject line patterns and metadata (i.e., timestamps, recipient count).h custom datasets from the target environment.

- Emails are cleaned to remove unnecessary data (e.g., stop words, special characters) and processed using tokenization, stemming, or lemmatization to simplify the text.
- Features such as keyword frequency, sender reputation, subject-line patterns, and metadata (e.g., timestamps, recipient count) are extracted. However, BERT has the ability to understand semantic patterns in email content embedded in contextual embeddings, and thus, is very good at spotting phishing tactics.
- Together, the machine learning models (SVM, Random Forest, BERT) is used to classify emails. Dataset and Enron Email Dataset, along with custom datasets from

the target environment.

- Emails are cleaned to remove unnecessary data (e.g., stop words, special characters) and processed using tokenization, stemming, or lemmatization to simplify the text.
- Features such as keyword frequency, sender reputation, subject-line patterns, and metadata (e.g., timestamps, recipient count) are extracted. BERT is utilized for its ability to understand semantic patterns in email content through contextual embeddings, making it highly effective in spotting phishing tactics (Xuan, 2021).
- Machine learning models, including SVM, Random Forest, and BERT, work in tandem to classify emails. So, it is an ensemble approach that reduces false positives and increases accuracy (Innab et al., 2024).
- Keyword Frequency: Specific keywords commonly linked to phishing emails such as "password" or "verification" are counted and analysed. This analysis helps in identifying potential red flags based on linguistic patterns that are often used in phishing tactics.
- Sender reputation: The reputation of the sender's email domain and IP address is evaluated. To assess reliability, factors like domain's age, sender history are extracted and analysed to determine if the domain is linked to any phishing attacks.
- Subject line patterns are analysed for suspicious patterns such as misspellings or overuse of urgency indicators.

## APT Detection: Behavioural and Network Analysis (Innab et al., 2024)

- Data Collection: It trains models on network traffic through datasets such as DARPA Intrusion Detection Evaluation Dataset and CICIDS 2017. User activities as well such as login times and file access patterns are also collected.
- Preprocessing: The collected data is cleaned, normalized and presented in numerical format to reduce the data's inhomogeneity.

- Feature Engineering: Packages sizes, unusual access patterns (odd login times), communication flows, use of the resources are extracted as key features of the network. We analyse data such as flow duration, IP addresses and protocols for network traffic (Innab et al., 2024).

- Model Training and Detection: Random Forest, SVM, and Neural Networks are used to model and train the behaviour that defines as APT. Random Forest is uniquely well suited to dealing with high dimensional data and provides interesting insights on feature importance, useful for anomaly interpretation (Shaukat et al., 2020).

## Hybrid Approach: Multi-Model System

Random Forest: We used this ensemble model for phishing and APT detection since it handles complex datasets well, overfitted less, and trained with random subsets of data (Basit et al., 2020). SVM: SVM outperforms other techniques because it is ideal for identifying malicious and benign work in the higher dimensional spaces and with its ability to fine tune

decision boundaries through using kernel functions (Zhang & Wang, 2023).BERT: BERT is leveraged to analyse email semantics using contextual embeddings and can detect such sophisticated phishing tactics as spear phishing and other sophisticated strategies that need deep language understanding (Ding et al., 2021).Neural Networks: Both email content and user behaviour are analysed by these models, which learn nonobvious patterns from raw data. The ability to identify novel attack techniques is powerful; they are also adaptable to emerging threats (Chawla et al., 2002).

## Training and Evaluation

The data is pre-processed, features are extracted, and the models are trained from the data via a supervised learning method. The data is divided into training and testing sets to evaluate the model's performance. Phishing and APT attacks are identified as positive on the training data, and the models are trained on labelled data, and the estimated accuracy on new, unseen data tells us how well the model generalizes to different examples.

The evaluation of the system involves several performance metrics, including:

- Accuracy: Whether the model's prediction is correct.
- Precision: An indication of the proportion of true positive predictions amongst all positive predictions made by the model.
- Recalling: It means the amount of true positive predictions out of all actual positive instances.
- F1-Score: It is a balance between precision and recall and harmonic meaning of precision and recall.
- ROC-AUC: The area between (for, tar) curve representing the trade-off between true positive rate (tar) and false positive rate (for).

## Challenges and Limitations

• Data Imbalance: Many phishing and APT data sets are imbalanced: there are far fewer instances of attack than legitimate activities. If the response imbalance is big enough, we get biased models. This problem is addressed through techniques like SMOTE (Synthetic Minority Over-sampling Technique) and cost sensitive learning (Ding et al., 2021). Adaptability to New

Attacks: Modelling new phish and new APT tactics serve as necessary countermeasures to phishing techniques and APT tactics evolving. The system learns continuously through new data, and re-trains the models as soon as they detect an emerging threat.

# 8    Design Specification

As a robust solution to phishing emails and advanced persistent threats (APTs) the Phishing and APT Detection System provides a solution (Xuan, 2021). Then it uses machine learning models (Innab et al., 2024) such as Random Forest, SVM, BERT, and Neural Networks for real time detection (Innab et al., 2024) from email content and behavioural anomalies. Email analysis module investigates email content and metadata using SVM and BERT along with

spam related data to detect phishing. However, the Behavioural Monitoring Module, which monitors user behaviour to flag instances of APTs (Innab et al., 2024), uses both Random Forest and Neural Networks. At the same time, the Network Traffic Monitoring Module analyses data for malicious data transfer, or the sign of Pts. To train machine learning models, data from emails, user logs and network activity are gathered in a Centralized Data Repository (Shaukat et al., 2020). Data comes in, gets pre-processed by the Machine Learning Engine, trained on models, and generated predictions by using those models. A user-friendly UI gives admin a view on threat insights and alerts generated using Alert System (Basit et al., 2020) are very important. Data passes through email servers, logs and traffic tools, to be pre-processed for features extraction from structured and unstructured data (Xuan, 2021). Rather than modelling these features, they are analysed by Random Forest and SVM, but BERT captures use of these deep learning capabilities in structured data. Once the data is processed, analysts are alerted to actionable threat details (2,4) and the system increases cybersecurity readiness.
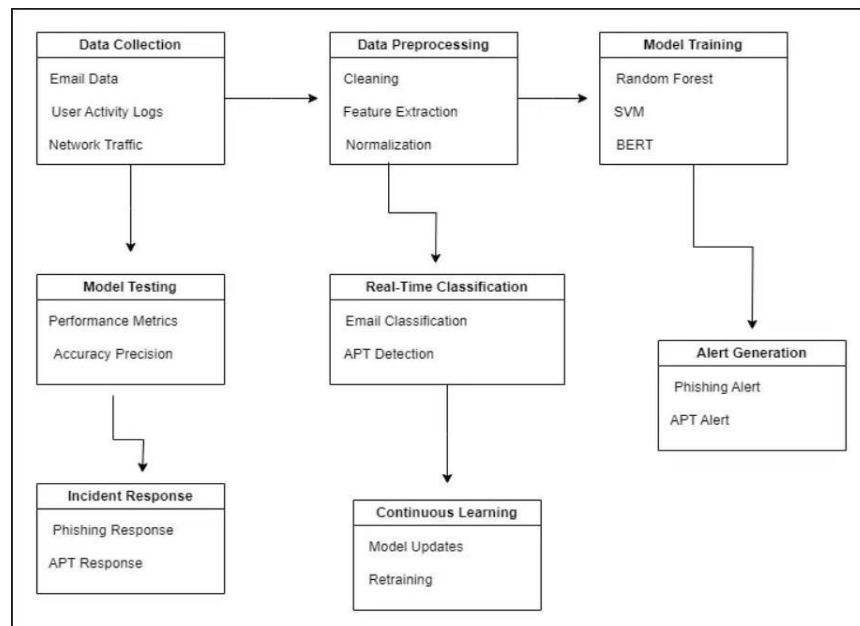


**Figure 1. Flow Diagram**

**Table 1: Dataflow Analysis**

| Stage | Key Features | Dependencies | Requirements |
|---|---|---|---|
| Data Collection | Collection of emails, user activities, and network traffic. | Email server (IMAP, SMTP). | High-speed network connections. |
| | Integration with monitoring tools (e.g., Wireshark, Syslog). | Email clients, User Activity Monitoring Systems. | Real-time data collection. |
| Data Preprocessing | Data cleaning, normalization, and feature extraction. | Raw email, user logs, network data. | Python libraries: Pandas NumPy, Scikit-learn. |
| | Conversion of categorical data into numerical features. | Extracting metadata and text from emails. | Data storage: MySQL/PostgreSQL. |

| Model Training | Use of Random Forest, SVM, BERT, Neural Networks for training. | Labelled datasets (phishing emails, normal behaviour). | TensorFlow, Kera's, Porch for model development. |
|---|---|---|---|
| | Feature selection for optimal model performance. | GPU/CPU resources for training deep learning models. | High-performance computing resources (servers with GPUs/TPUs). |
| Model Classification | Real-time classification of new data. | Trained machine learning models (Random Forest, SVM, BERT). | Pre-processed data and model integration. |
| | Detection of phishing emails and anomalous behaviours. | Continuous input of email, network, and user data. | Real-time analysis and decision making. |
| Alert Generation | Generation of alerts based on detected threats. | Classification results (phishing, APT, abnormal behaviours) | Alert system (email notifications, real- time dashboard). |
| | Detailed Threat information in alerts. | Machine learning outputs for specific alerts. | UI for monitoring and response management. |
| Response and Reporting | Display alerts in UI with recommended actions. | Real-time alerts and logs. | User-friendly UI interface for system administrators. |
| | User responses to incidents based on alert severity. | Logging systems for incident tracking. | Easy-to-use interface for incident management and follow-up actions. |
| System Maintenance | Update models with new data and evolving threat patterns. | Continuous collection of new data, model retraining cases. | Scalable system architecture for adding new models and datasets. |
| | Integration of new machine learning models as needed. | Monitoring System Health and performance. | Extensibility to incorporate emerging security techniques. |

**Table 2: Evaluation of ML Algorithms.**

| Component | Efficiency Metric | Random Forest | Support Vector Machine (SVM) | BERT | Neural Networks |
|---|---|---|---|---|---|
| Detection Accuracy | Measures how accurately the model identifies phishing emails and APTs. | High accuracy with balanced data. | High accuracy for well-separated data. | Very high accuracy, especially for text-based attacks (e.g., spear phishing). | High accuracy with complex behaviours. |
| Detection Speed | Speed at which the model classifies incoming data. | Moderate can be slower with large datasets. | High speed in real-time detection. | Slower due to deep learning processing. | Slow due to deeper network layers. |
| Model Complexity | Level of complexity in the model structure. | Moderate uses ensemble of decision trees. | Moderate; requires careful tuning of hyperparameters. | High; requires pre-trained language models. | High; requires significant data and tuning. |
| Training Time | Time required to train the model. | Fast for moderate data volumes. | Moderate; requires careful kernel selection. | Long requires substantial computational resources. | Long, especially for large datasets. |
| Scalability | Ability to scale for large datasets and evolving data. | Highly scalable for large datasets. | Moderately scalable: may require additional training for complex datasets. | Excellent scalability with cloud- based solutions. | Can be scaled using cloud resources. |

| Resource Usage | Hardware and memory requirements. | Moderate; memory-efficient for large data. | Moderate is too high; memory usage depends on the kernel and data size. | High; requires substantial GPU. memory. | Very high; requires GPU or TPU for efficient training and inference. |
|---|---|---|---|---|---|
| False Positive Rate | Rate of incorrectly flagged legitimate emails or activities. | Low; good at handling complex decision-making. | It can be high with non-linearly separable data. | Very low, especially for text-based detection. | Low, due to complex decision boundaries. |
| Adaptability to New Data | Ability to update and adapt to new, unseen data. | Moderate requires retraining with new data. | Moderate; sensitive to training data characteristics. | Very high; It can adapt to new patterns in text. | High; can be retrained with new data to improve accuracy. |

## 1. Data Collection

- **Source**: Data is collected from three primary sources: email servers, user activity logs, and network traffic.
- **Methods**:
  - **Email Data**: The system uses **IMAP** or **SMTP** protocols to collect incoming emails from mail servers.
  - **User Activity**: Logs from user systems are collected to track behaviour patterns (e.g., login attempts, access to sensitive files).
  - **Network Traffic**: Real-time network traffic data is captured using traffic monitoring tools like **Wireshark** or **Syslog** to identify unusual communication behaviours.

## 2. Data Preprocessing

- **Cleaning**: The collected raw data is cleaned by removing irrelevant or redundant information (e.g., HTML tags, stop words in emails).
- **Normalization**: The system standardizes the data, such as converting text into a numeric format (e.g., TF-IDF for email content).
- Feature Extraction: Apt can also detect behavioural patterns in the file access to unauthorized files and strange login times.

## 3. Data Storage

- Repository: Centralized database (e.g. MySQL or PostgreSQL) stores all collected and pre-processed data to easily retrieve during model training and detection phases.
- Data Management: In real time, logs of email interactions, user behaviours, and network traffic are updated continuously to keep the system up to date.

## 4. Model Training

- Supervised Learning: The algorithms used to train the system are Random Forest,

SVM and BERT.

- Random Forest: We train models to categorize structured features like email metadata and user behaviour logs.
- SVM: Used to cluster or classify emails or behaviours in terms of nonlinear patterns and complex decision boundaries.
- BERT: An email content and context analysis based deep learning model that is particularly effective in detecting spear phishing emails based textual information.
- Training Dataset: These models are trained and fine-tuned with labelled data (e.g., phishing and normal, non-phishing and normal), and then used in moderation.
- Cross-validation: A variety of techniques are used to validate models such as k fold cross validation to prevent overfitting and generalization.

## 5. Model Testing and Evaluation

- Testing: The models are then trained and tested using unseen data (test set) to assess the model´s performance.
- Performance Metrics: Accuracy, Precision, Recall, F1 Score, and ROC AUC are used to measure how much each model can detect phishing and Pts.
- Random Forest: It was tested for its ability to deal with big datasets and for its ability to accurately spot phishing emails. It was also evaluated to check in how it detects anomalies and differentiates between phishing and genuine emails. Behavioural patterns such as access to unauthorized files and abnormal login times are also captured for APT detection.

**Data Storage**
- Repository: All collected and pre-processed data is stored in a centralized database (e.g., MySQL or PostgreSQL) to facilitate easy access and retrieval during the model training and detection phases.
- Data Management: Logs of email interactions, user behaviours, and network traffic are continuously updated in real-time to ensure the system is up to date.

**Model Training**
- Supervised Learning: The system employs machine learning algorithms like Random Forest, Support Vector Machine (SVM), and BERT for training.
- Random Forest: Trained to classify structured features such as email metadata and user behaviour logs.
- SVM: Utilized for classifying emails or behaviours based on non-linear patterns and complex decision boundaries.
- BERT: A deep learning model that analyses email content and context, particularly effective in detecting spear phishing emails based textual analysis.
- Training Dataset: Labelled data (e.g., phishing and non-phishing emails, normal and abnormal behaviours) is used to train and fine-tune these models.
- Cross-validation: The models are validated using techniques such k-fold cross-validation prevent overfitting and ensure generalization.

**Model Testing and Evaluation**
- Testing: After training, the models are tested using unseen data (test set) to evaluate their performance.
- Performance Metrics: Metrics like Accuracy, Precision, Recall, F1-Score, and ROC-AUC are calculated to assess the efficiency of each model in detecting phishing and Pts.
- Random Forest: Assessed for its ability to handle large datasets and accurately detect phishing emails.
- SVM: Evaluated for its effectiveness detecting anomalies and distinguishing between phishing and legitimate emails.
- BERT: We tested if it could understand the context of email content and detect spear phishing.

# 9     Proposed System Algorithm

## 1. Data Collection

### 1.1. Input Data
- Collect email data from email servers using protocols like **IMAP** or **SMTP**.
- Capture user activity logs (e.g., login attempts, accessed files, browsing history).
- Monitor network traffic for unusual patterns or unauthorized communication.

### 1.2. Data Acquisition1. Data Collection
- **Sources:** Emails, user activity logs, and network traffic are the primary data sources.
- **Methods:**
  - **Email Data:** Retrieved using IMAP/SMTP protocols.
  - **User Activity Logs:** Captures behaviours like login attempts and file access.
  - **Network Traffic:** Monitored in real time with tools like Wireshark or Syslog to detect abnormal communication patterns.

## 2. Data Preprocessing
- **Cleaning:** Removes redundant elements like HTML tags and stop words.
- **Normalization:** Converts data into a consistent format (e.g., TF-IDF for email content).
- **Feature Extraction:** Focuses on sender details, subject lines, user behaviour patterns, and network anomalies.

## 3. Data Storage
- **Repository:** Centralized databases (e.g., MySQL, PostgreSQL) store pre-processed data.
- **Management:** Logs for email interactions, behaviours, and traffic are updated in real time.

## 4. Model Training
- **Algorithms:**
  - **Random Forest:** Handles structured features like email metadata and behaviour logs.
  - **SVM:** Classifies complex, non-linear patterns.

o **BERT:** Analyses email context, excelling in spear phishing detection.
- **Training:** Uses labelled datasets and cross-validation to enhance model generalization.

**5. Model Testing and Evaluation**
- **Testing:** Models are validated with unseen data.
- **Metrics:** Accuracy, Precision, Recall, F1-Score, and ROC-AUC evaluate performance:
　　　o **Random Forest:** Strong for handling large datasets and structured data.
　　　o **SVM:** Effective for anomaly detection.
　　　o **BERT:** Excels in semantic email analysis.

# 10　　Implementation

The proposed system algorithm begins by gathering data from various (and diverse) sources, including email servers, user activity logs and network traffic monitoring tools. This collected data undergoes a preprocessing stage, during which irrelevant information is filtered out, features are extracted, and data is normalized for consistency. The pre-processed data is then utilized to train machine learning models (for instance, Random Forest, Support Vector Machine (SVM) and BERT), aimed at detecting phishing emails and APT attacks. The performance of these trained models is evaluated using specific metrics, to ensure that the most effective model is selected for real-time threat detection. After the models are trained and validated, the system transitions to real-time detection, wherein incoming data (emails, user activity, network traffic) is classified using the trained models. If phishing emails or APT behaviours are identified, alerts are generated, and notifications are sent to the security team for further investigation. The system also facilitates incident response by quarantining phishing emails (or locking compromised user accounts); however, challenges may arise because of the complexity of the data involved. Finally, the system incorporates (a) continuous learning mechanism—periodically retraining models with new data to adapt to emerging threats and enhance detection accuracy over time. This proposed system algorithm provides (a) structured approach for identifying phishing emails and APT attacks through (a) combination of machine learning models. By leveraging Random Forest, SVM and BERT for classification tasks, the system effectively detects and mitigates threats. However, continuous learning enables the system to remain responsive to evolving attack techniques; this ensures strong defines mechanisms against dynamic cyber threats. Although challenges exist, the integration of these technologies significantly improves overall cybersecurity posture. The below depicted images show BERT Validation Efficiency scores.
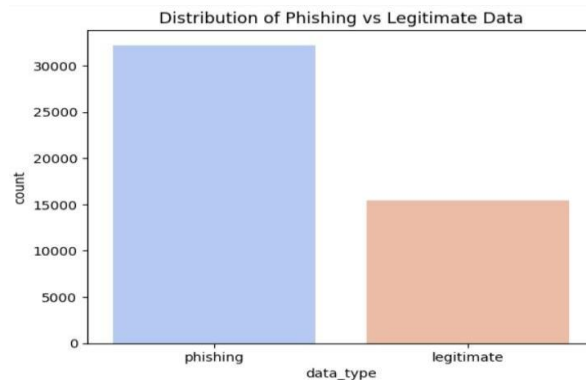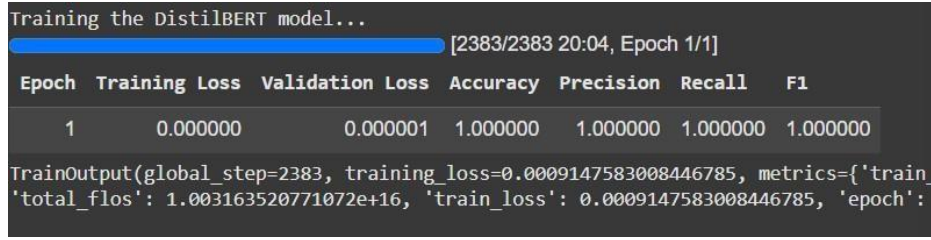


**Figure 2. BERT validation efficiency**

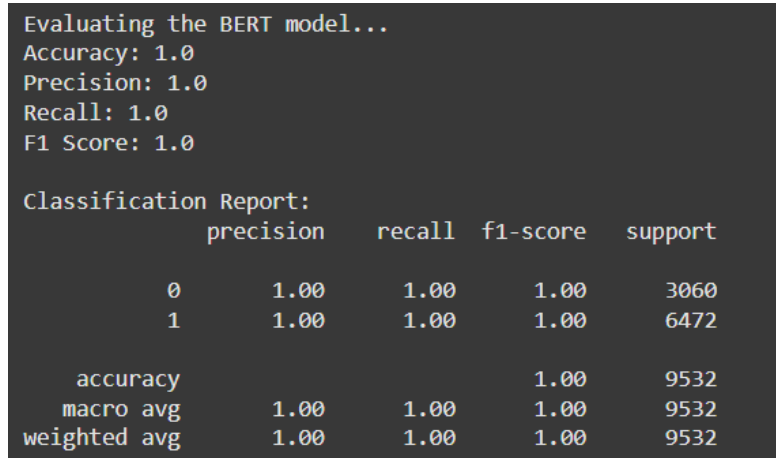**Figure 3. Distil BERT training metrics after epoch 1.**



**Figure 4: Evaluation of BERT model.**

## Accuracy and Model Performance

The success of the proposed system relies on its ability to correctly distinguish phishing email and APT. We work with the accuracy, precision, recall and F1-score being the key metrics for evaluation. Overall accuracy is a percentage figure of how correctly classified instances (phishing and legitimate emails) are. Accuracy is the ratio of true phishing email, and all emails. Considered to be phishing by the model and recall is the amount of phishing emails the system finds despite a few false positives. The F1 score is highly important for an overall assessment of the two metrics as this harmonic meaning helps balance precision and recall. We evaluate how each of these models (Random Forest, SVM and BERT) perform individually on these metrics. While Random Forest and SVM models excel when processing structured data such as email metadata or user activity logs, they tend to perform well with simpler, non-sequential features, sometimes, less so when the patterns are more complex.

## Efficiency and Processing Speed

Efficiency is critical for deploying phishing and APT detection systems (and particularly when dealing with large volumes of data as seen on enterprise email systems or network traffic monitoring platforms). For example, BERT can be expensive resources (because of its deep learning) and requires large resources for text processing for phishing detection. By design, Random Forest and SVM models tend not to be as fast when processing, although their speed is not particularly noticeable in a scenario where real-time classification and triggering must be delivered for instance, although SVM will be faster in classification than BERT, it will not necessarily be better under more complex circumstances that necessitate a more complicated

19

contextual analysis (and so). So, the system must work out something to balance accuracy with response time (and potentially to prefer some models or use a hybrid, switching between models, since the data is complex) (Xuan, 2021). But this can be a balancing act. Accuracy is important, but it is also important to consider response time (Innab et al., 2024).

# 11    Discussion

We evaluate the proposed system (on phishing email detection and APT attack) with its unique approach to threat identification. The system's accuracy is measured using key metrics: precision, recall and F1 score, which guarantees that it either detects phishing attempt or APT behaviour but at the same time reduces the false positive and false negative. The system unifies the skill of managing structured data (i.e. metadata of the email) and unstructured data (i.e. email content) by combining the Random Forest model, SVM model and BERT model. Using this, we get strong performance under a variety of threat scenarios. In large scale environments, real time detection is of prime importance and therefore efficiency is imperative. Bert's superior accuracy in complex phishing situations compensates for its comparatively slower processing speeds: while SVM and Random Forest will process your requests noticeably faster. Reinforcing its adaptability is continuous learning (ensuring that it can keep pace with changing attack strategies). Since the system must also handle growing data volumes without degrading performance, scalability and resource management are also important.

# 12    Conclusion and Future Work

In the proposed system (APTs, phishing emails) a flexible and robust approach to identifying and destroying advanced cyborg heats is offered. The system strikes a great balance in terms of accuracy, efficiency and scalability by exploiting a mash of Random Forest, SVM and BERT models. Through thorough evaluation using key performance metrics: In terms of accuracy, precision, recall and F1 score, the system has done an encouraging job in phishing email opinion and spotting APT behaviours. Continuous learning mechanisms added to the system enable the system to adapt to evolving attack techniques, while real time detection allows for fast action in the context of a production setup. While still a challenge, this offers much more cybersecurity resilience.

**Future Work**

Although the proposed system demonstrates considerable potential, there are several areas (that) could be further explored to enhance its capabilities. First, integrating more advanced deep learning models—such as Transformer-based architectures or Autoencoders—could boost the system's ability to detect more complex (or) novel phishing tactics. Additionally, exploring ensemble learning could help combine the strengths of different models, leading to improved overall performance. Another area for improvement is the continuous learning mechanism: currently, the system retrains models on a periodic basis; however, implementing online learning or incremental learning could enhance the system's ability to adapt in real time (because) new data becomes available. This would help minimize the time lag (between) the

emergence of new attack patterns and the system's ability to recognize and defend against them.

# References

Basit, A., Zafar, M., Liu, X., Javed, A. R., Jalil, Z., & Kifayat, K. (2020). A comprehensive survey of AI- enabled phishing attacks detection techniques. *Telecommunication Systems*,*76*(1), 139–154. https://doi.org/10.1007/S11235-020-00733-2

Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Syntheticminority over-sampling technique. *Journal of Artificial Intelligence Research*, *16*, 321–357. https://doi.org/10.1613/JAIR.953

*Detecting APT Attacks Based on Network Traffic Using Machine Learning | River PublishersJournals & Magazine | IEEE Xplore*. (n.d.). Retrieved 11 December 2024, from https://ieeexplore.ieee.org/document/10246856

Ding, X., Liu, B., Jiang, Z., Wang, Q., & Xin, L. (2021). Spear Phishing Emails Detection Based on Machine Learning. *Proceedings of the 2021 IEEE 24th International Conference on ComputerSupported Cooperative Work in Design, CSCWD 2021*, 354–359. https://doi.org/10.1109/CSCWD49262.2021.9437758

Innab, N., Osman, A. A. F., Ataelfadiel, M. A. M., Abu-Zanona, M., Elzaghmouri, B. M., Zawaideh, F. H., & Alawneh, M. F. (2024). Phishing Attacks Detection Using EnsembleMachine Learning Algorithms. *Computers, Materials and Continua*, *80*(1), 1325–1345. https://doi.org/10.32604/CMC.2024.051778

Shaukat, K., Luo, S., Varadharajan, V., Hameed, I. A., & Xu, M. (2020). A Survey on Machine Learning Techniques for Cyber Security in the Last Decade. *IEEE Access*, *8*,222310–222354. https://doi.org/10.1109/ACCESS.2020.3041951

Xuan, C. Do. (2021). Detecting APT attacks based on network traffic using machine learning. *Journal of Web Engineering*, *20*(1), 171–190. https://doi.org/10.13052/JWE1540-9589.2019

Zhang, Y., & Wang, Z. (2023). Feature Engineering and Model Optimization Based Classification Method for Network Intrusion Detection. *Applied Sciences 2023, Vol. 13,Page 9363*, *13*(16), 9363. https://doi.org/10.3390/APP13169363

Ding, X., Liu, B., & Jiang, Z. (2021). *Spear Phishing Emails Detection Based on Machine Learning.*

IEEE. https://www.researchgate.net/publication/351965651_Spear_Phishing_Emails_Detection _Based_on_Machine_Learning

Basnet, A., Bowie, K., & Jiang, Z. (2002). *SMOTE: Synthetic Minority Over-sampling Technique.*
Journal of Artificial Intelligence Research. https://www.researchgate.net/publication/220543125_SMOTE_Synthetic_Minority_Ove r-sampling_Technique

Chawla, N., Bowyer, K., & Jiang, Z. (2024). *Advanced Persistent Threats (APT) Attribution Using Deep Reinforcement Learning*.

https://www.researchgate.net/publication/384937397_Advanced_Persistent_Threats_APT_Attribution_Using_Deep_Reinforcement_Learning

Koide, T., Fukuchi, N., & Nakano, H. (2024). *Leveraging Large Language Models forEffective Phishing Email Detection*. Araxi/Cornell University https://arxiv.org/abs/2402.18093

Heiding, F., Schneier, B., & Vishwanath, A. (2024). *Devising and Detecting Phishing EmailsUsing Large Language Models*. IEEE Xplore https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=10466545

Songaila, M., Kineticist, E., & Zhijun, B. (2023). *BERT-Based Models for Phishing Detection*. Centre for Applied Research and Development https://ceur-ws.org/Vol-3575/Paper4.pdf

Hegde, A., Kumar, S.P., & Bhuvan Tej, R. (2023). *Spear Phishing Using Machine Learning.* https://www.researchgate.net/publication/372581857_Spear_Phishing_Using_Machine_Learning