# MACHINE LEARNING AND PORT SCANS: A SYSTEMATIC REVIEW

**Jason M. Pittman**
ORCID: 0000-0002-5198-8157

## ABSTRACT

Port scanning is the process of attempting to connect to various network ports on a computing endpoint to determine which ports are open and which services are running on them. It is a common method used by hackers to identify vulnerabilities in a network or system. By determining which ports are open, an attacker can identify which services and applications are running on a device and potentially exploit any known vulnerabilities in those services. Consequently, it is important to detect port scanning because it is often the first step in a cyber attack. By identifying port scanning attempts, cybersecurity professionals can take proactive measures to protect the systems and networks before an attacker has a chance to exploit any vulnerabilities. Against this background, researchers have worked for over a decade to develop robust methods to detect port scanning. While there have been various surveys, none have focused solely on machine learning based detection schemes specific to port scans. Accordingly, we provide a systematic review of 15 papers published between February 2021 and January 2023. We extract critical information such as training dataset, algorithm used, technique, and model accuracy. We also collect unresolved challenges and ideas for future work. The outcomes are significant for researchers looking to step off from the latest work and for practitioners interested in novel mechanisms to detect the early stages of cyber attack.

*Keywords*  systematic review, machine learning, port scanning, cybersecurity, algorithms, training data

## 1 Introduction

Cybersecurity incidents continue to plague digital life. While a significant portion of incidents result from phishing and malware, 45% are the result of network-based cyber attacks [1]. These cyber attacks follow a pattern or procedure. Existing models and methodologies vary in the number of steps. However, the first step is universally understood to be *reconnaissance*. In turn, reconnaissance most often includes some type of port scanning.

Port scanning is a technique to enumerate target endpoints. Confusingly, port scanning can be both a legitimate engagement [2] or a malicious precursor to escalating intrusion [3, 4]. A general issue is differentiating between what may be an authorized benign instance of host enumeration and a malicious scanning of active hosts and their available ports. Furthermore, if we accept port scanning as a necessary prelude to cyber attack, then we want to develop a means to detect port scanning with high certainty. To this end, there is a small but growing literature on detecting port scanning. The literature ranges from early intrusion detection mechanisms [5] to sophisticated machine

learning techniques [6]. There have been several comparative surveys during this time, most recently Aamir et al. [7] and [8]. However, there has not been a systematic review of the literature.

Literature reviews are invaluable to a field of study. Reviews provide an understanding of the existing research by establishing a foundation of knowledge. Reviews also clarify existing knowledge related to a given problem. Both functions guide new investigations and reduce overlap or unnecessary duplication of work. Yet, new reviews are necessary as the field grows, new techniques are discovered, and new technologies are released which impact forms of inquiry. Accordingly, the purpose of this study is to provide a systematic review of existing literature using machine learning algorithms to detect port scanning.

The remainder of this work is organized in a way which (a) situates the systematic review in existing knowledge and (b) maximizes understanding of the cutting edge. The first is achieved by discussing port scanning and detection of such port scanning literature. Thereafter, we present the research method and techniques used to find, organize, and analyze research published since 2021. Finally, we

demonstrate the findings of the analysis in terms of quantitative results from the existing research.

## 2 Related Work

The work most proximal to this study exists in two categories: scanning TCP/IP ports and detection of those scans. The following discussion is not intended to be exhaustive. Rather, we offer background research that we view as seminal and salient.

### 2.1 Port Scanning

Port scanning uses features of TCP/IP to enumerate computing systems on a network. As different network protocols use different ports, it's essential to scan a wide range of ports to gather complete information. This is because vulnerabilities can exist in all protocols. The total number of ports that can be scanned is 65535, with ports 0 to 1023 being well-known, ports 1024 to 49151 being registered, and ports 49152 to 65535 being dynamic or private.

The origin of the phrase *port scanning* in the academic literature can be traced back to the early days of computer networking. In the late 1980s and early 1990s, as the Internet was growing and becoming more widely used, there was an increasing need for tools to help network administrators and security professionals understand the state of their networks. One of the key tasks for these professionals was identifying which network services were running on which hosts, and which ports on those hosts were open or closed. This process became known as *port scanning*.

One of the earliest references to port scanning in the literature is found in Fyodor [9], which described a method for determining the operating system of a remote host by sending probes to specific ports and analyzing the responses. The work explained the operating system of a host can be determined by analyzing the TCP/IP stack's behavior and its responses to different types of probes, such as the initial sequence numbers (ISNs) and the options in the TCP headers of the responses. Further, the paper also described how to use this technique to fingerprint the operating system of a remote host as well as the limitations and challenges of the technique. Additionally, the paper introduced the first version of an open-source tool named *nmap* (Network Mapper) that implements this technique for Remote OS detection. While nmap is not the only port scanner available, it is featured heavily throughout the literature.

De Vivo et al. [10] generalizes from the port scanning foundation provided in Fyodor [9] and several [11, 12] others. The significance of De Vivo et al. [10] emerges from the rigorous classification applied to port scanning techniques and procedures. The paper described the different types of port scans, such as TCP connect scans and SYN scans as *classical*. This is in relation to *indirect* and *stealth* scanning. The latter is also referred to as a FIN, XMAS, or NULL scan. The former is realized by bounc-

ing scans off of a zombie endpoint. The work goes on to describe scanning techniques. These includes *decoy* scanning, *fragmented* scanning, and *coordinated* or *distributed* scanning, UDP scanning, and ICMP sweeping.

Barnett et al. [13] presented a classification system for network scanning techniques. The significance of the work is in establishing a clear and organized classification of the different types of network scanning techniques that exist and their use cases. To that end, the authors propose a taxonomy categorizing network scanning techniques based on the level of interaction with the target system, the type of information gathered, and the purpose of the scan. This extends De Vivo et al. [10] in both types and techniques.

Barnett et al. [13] add two additional network scanning techniques to the three presented by De Vivo et al. These are *vertical*, *horizontal*. Further, Barnett et al. differentiate between OSI Model layer 2 scans and layer 3 scans with overlaying attributes according to speed (slow, medium, rapid) and distribution (one-to-one, one-to-many, many-to-one, many-to-many). Barnett et al. also describe how scan types from prior work (e.g., De Vivo et al.) map to their categories. The mapping is more pronounced when attributes encompassing speed and distribution were considered.

Bou et al. [14] demonstrated a comprehensive overview of the different types of cyber scanning techniques that are used to identify various features of networks. The authors divide port scanning techniques into two main categories: *passive* and *active*. Passive scanning techniques involve listening to network traffic to gather information about the target network without sending any packets. Active scanning techniques involve sending packets to a target host to elicit a response, which can be used to determine the host's characteristics and identify vulnerabilities.

The work extends the categorization by describing different techniques of passive and active scanning techniques. This calls back to the organizational structure provided by De Vivo et al. [10] and Barnett et al. [13] but differs in semantics. For instance, the De Vivo et al. classical, indirect, and stealth scans map under the nature of active and passive scanning. Further, the semantic developed by [13] around relations between scanner and target (e.g., one-to-many) falls under *approach* in Bou et al. [14]. Bou et al. also offered *strategy* as a way to categorize directional relationship between scanner and target.

Roy et al. [15] claimed a gap exists in the literature on classifying and categorizing adversarial reconnaissance processes. The claim stands to reason given the authors first delineate between *technical* and *non-technical* reconnaissance techniques. Technical reconnaissance included network scanning, or *cyber scanning* as Roy et al. refer to it, as a remote technique. The authors then differentiate between *host detection* and *port enumeration* which stand as a combined label for all of the scanning techniques outlined by De Vivo et al [10] (i.e., ICMP, SYN, Full Connect, etc.).

2

While Roy et al. [15] did not add anything new to the port scanning taxonomy, the authors did connect the prior research by De Vivo et al. [10] and Barnett et al. [13] to a burgeoning literature around detection of port scanning.

## 2.2 Detecting Port Scanning

Port scanning, as a reconnaissance technique, is detectable. ML is a compelling solution to detecting otherwise undetectable port scans because of its ability to correlate seemingly unrelated features across enormous datasets. Yet, not all ML algorithms work in the same way or have the capability to address the same problem. Furthermore, there are a variety of ML algorithms types-classification, regression, deep learning, and so forth- with a diversity of implementation variations.

The majority of works investigating ML for detecting port scanning over the past decade and a half include at a review of prior ML algorithm performance. Such work exists as a quasi review with an add-on quantitative analysis of an algorithm not featured in the prior research. We use the term *quasi* because these works review algorithms by running each against a common training and evaluation dataset. These studies do not rely on results from the prior research responsible for introducing the algorithm to the field. This is close to the notion of a meta-analysis but not precisely so. Still, the quasi reviews are particularly significant for researchers and practitioners looking to get up to speed on the state of the field in short order. We discovered two such works and include those as foundational literature which we directly extend with our systematic review.

Aamir et al. [7] investigated the detection of characteristics of port scanning and analyzed the performance of 22 ML algorithms. The algorithms included decision trees, discriminant analysis, support vector machines (SVM), k-nearest neighbors (KNN), and ensemble classifiers. The authors used the CICIDS2017 dataset with a 70%training and 30% evaluation split. Of the 22 ML algorithms examined, nine demonstrated more than 85% classification (testing) accuracy. Specifically, Aamir et al. identified Fine Gaussian SVM as best performing algorithm with 99% testing and 99% training accuracy. False negative rates are provided for all 22 algorithm experiments. Further, for fast training with high accuracy scores, discriminant analysis was more accurate and efficient in classifying port scans. Aamir et al. did not discuss the type of port scans detected nor what scanning techniques were present in the dataset.

Krishna et al. [8] also investigated port scan detection using a variety of ML algorithms. The authors analyzed fewer algorithms but used the same training and evaluation dataset as Aamir et al. [7]. Krishna et al. examined two of the algorithms as Aamir et al., SVM and decision trees. Krishna et al. also evaluated random forest and logistic regression. Unlike any other study we found in the literature, Krishna et al. do not present results. Instead,

the authors include snapshots of their Juypter notebook as figures. On one hand, including code makes the work repeatable and reproducible. On the other hand, one would need to repeat and reproduce the study to obtain algorithm performance values.

Both Aamir et al. [7] and Krishna et al. [8] represent the predominant type of research in the literature. That is, existing port scan detection research frequently examines ML algorithm performance by direct experimentation rather than reference to prior experimentation. Consequently, the findings from these studies are scattered throughout the literature. Anyone interested in extending the field is left to trace through the forest to find relevant trees. With that in mind, the goal of this work was to systematically review results published since 2021 to catalog ML algorithm performance in detecting port scans.

## 3 Method

This work employed a systematic literature review methodology. Systematic literature review is a well-defined method that is used to identify, evaluate, and interpret all of the available research on a particular topic [16, 17]. The process is designed to be comprehensive, unbiased, and transparent, and it involves a number of steps, including formulating a research question, searching for relevant literature, selecting studies, extracting data, and synthesizing the results [16, 18]. Systematic literature reviews are increasingly used in the field of software engineering and other technical fields, but also in other scientific fields, as a way to provide an in-depth understanding of existing knowledge on a topic.

A systematic literature review differs from other literature review methods in several ways. A traditional literature review, also called a narrative review [16], is typically less structured and less systematic. It is often used to provide an overview of the current state of knowledge on a topic, but it is not as rigorous as an SLR in terms of the search and selection process.

With the design of systematic reviews in mind, we pose four questions.

RQ1: What machine learning algorithms have been used to detect port scanning?

RQ2: What were the detection rates and false positive rates for those algorithms?

RQ3: What datasets were used for training and evaluation of those algorithms?

RQ4: What port scanning types and techniques were used for evaluation of those algorithms?

We constrained the literature search to 2021 and newer. We did so based on the last relevant reviews being published in 2021 [7, 8]. While we recognize the majority of work in detecting port scans exists between 2010 and 2021, research is still progressing in this research area

and a systematic review of published research since 2021 holds significance for researchers and practitioners alike.

No literature search will produce perfect results. However, careful attention to search strings can yield sufficient results so as to be thorough. We used *"detecting port scan" AND "machine learning"* as a starting search string. The search returned 61 papers. For comparison, searching with *"port scan" AND "machine learning"* produced 952 results. Meanwhile, *"detect port scan" AND "machine learning"* produced 21 results. Often researchers use the commonly accepted short form of machine learning, ML. Thus, we used *"ML" AND "detecting port scan"* to cross-check the search. This produced 43 results. The variant of *"detect port scan" AND "ML"* found 12 papers.

We manually reviewed each work to ensure each study included detection of port scanning. Manual review was necessary because some work folds port scan detection into an overarching intrusion detection framework. We were left with 15 studies as our dataset after this step.

Data extraction from the selected papers is an important step to properly answer the research questions. In this study, we used the following data form to extract the needed information: (a) year of publication; (b) authors; (c) source of publication; (d) citation count; results (accuracy as F1); (e) dataset source; and (f) algorithms as task, technique, and procedure (TTP). We also included the port scanning types and techniques when such were available in the research.

## 4 Results

We separate the results of the systematic review into two sections. The first section provides an overview of our dataset. We describe the literature features for ease of future reference. Then, we present a breakdown of that literature by algorithm. As with Aamir et al. [7] and Krishna et al. [8], some work experimented with more than one ML algorithm. Such research appears in multiple categories below.

### 4.1 The Literature

We analyzed 15 studies published since 2021 (Table 1). Seven studies were from 2021, six were from 2022, and a single study appeared in early 2023. The remaining study was from 2020 which we included as a specific exception. This is discussed in the Neural Network (NN) algorithm section below.

The sample encompassed six total ML algorithms. The majority (10) of studies examined a single ML algorithm while five studies examined more than one. The literature published in 2021 spanned all six algorithms whereas literature from 2022 focused on a single algorithm (with one exception). Random Forest (RF) and SVM were the most investigated algorithms in the 2021 subset. A variety of NN implementations appeared throughout the 2022

subset. Nature-inspired (NI) appeared once while Regression (R) and Naive Bayes (NB) were studied three and five times respectively.

Six studies have not been cited. Six studies have been cited more than once with 25 being the highest citation count. Only one study [19] included a paper [7] from the literature population in its related work. The other 14 papers exist independent of one another with only indirect relations from support research in general ML or cybersecurity.

Table 1: Literature using ML to detect port scans

| Authors | Year | Cited | TTP |
|---|---|---|---|
| Hartpence et al. | 2020 | 6 | NN |
| Algaolahi et al. | 2021 | 1 | RF,SVM |
| Baah et al. | 2021 | 0 | RF,SVM,NB |
| Sirisha et al. | 2021 | 4 | RF,R,NB |
| Liu et al. | 2021 | 21 | NI |
| Bertoli et al. | 2021 | 25 | RF,SVM,R,NB,NN |
| Mohseni et al. | 2021 | 1 | RF |
| Al-Haija et al. | 2021 | 9 | NB |
| Bakaletz | 2022 | 0 | NB |
| Tojeiro et al. | 2022 | 0 | R |
| Singh et al. | 2022 | 2 | NN |
| Lv et al. | 2022 | 0 | NN |
| Kirtas et al. | 2022 | 1 | NN |
| SaiKiran et al. | 2022 | 0 | RF,SVM,NN |
| Henry et al. | 2023 | 0 | NN |

### 4.2 The Algorithms

We found six machine learning algorithms in the literature sample. The following sections present a summary for each algorithm and the relative meaning of using it to detect port scanning. We summarize each algorithms's performance in terms of accuracy and false positives. We also present the dataset used to train and evaluate the models when such are revealed in the source literature.

#### 4.2.1 Random Forest

Random Forest is good for classification problems, particularly in cases where there are many features and interactions among features. The algorithms is also useful for feature selection and handles missing data well.

Six studies out of the 15 study sample experimented with the RF algorithm [20, 21, 22, 23, 24, 25]. Algorithm performance ranged from 78.09% to 100% across those studies. One paper [21] included source code or a link to a source code repository (e.g., GitHub). Four different datasets were used, three of which do not appear in other algorithm categories.

Two studies discussed the types of port scans present in training and evaluation data. SaiKiran et al. [25] mentioned *port sweep* but did not specify further. Bertoli et al. [21] conducted training and evaluating against the full

spectrum of port scan types. Further, the authors included port scan data from five different port scan tools.

Table 2: Random Forest Algorithm Performance

| Authors | Accuracy (F1) | Dataset |
| --- | --- | --- |
| Algaolahi et al. | 99.75 | CICIDS2017 |
| Baah et al. | 99.98 | CICIDS2017 |
| Sirisha et al. | 78.09 | NSLKDD |
| Sirisha et al. | 84.14 | CICIDS2017 |
| SaiKiran et al. | 99.93 | CICIDS2017 |
| Mohseni et al. | 99.94 | CICIDS2017 |
| Bertoli et al. | 96.00 | MAWILab |
| Bertoli et al. | 100.00 | Bonafide |

### 4.2.2  Support Vector Machine (SVM)

Support Vector Machine (SVM) is good for classification and regression problems, especially in cases where the data has clear boundaries and is not noisy. It works well for datasets with a limited number of features.

SVM is different from Random Forest in that it uses a boundary (a hyperplane) to separate the data into classes, whereas Random Forest creates multiple decision trees and aggregates their predictions to make a final decision. SVM is best suited for cases where the boundary between classes is well defined and clear, whereas Random Forest is better suited for complex, non-linear decision boundaries.

Four studies experimented with SVM [20, 21, 24, 25]. All four also had explored RF performance. Results for the SVM experiments ranged from 89.61% to 99.87% both coming from the same dataset (of two total). Source code availability and port scan details remained the same as indicated in the RF algorithm category.

Table 3: Support Vector Machine Algorithm Performance

| Authors | Accuracy (F1) | Dataset |
| --- | --- | --- |
| Algaolahi et al. | 89.61 | CICIDS2017 |
| Baah et al | 99.87 | CICIDS2017 |
| SaiKiran et al. | 93.29 | CICIDS2017 |
| Bertoli et al. | 92.00 | Bonafide |

### 4.2.3  Regression

Regression algorithms are used for predicting a continuous target variable based on one or more input features. They are commonly used for tasks such as predictions and forecasting.

Regression algorithms, including linear regression, are different from Random Forest and SVM algorithms in that they focus on establishing a linear or non-linear relationship between the input features and the target variable. On the other hand, Random Forest and SVM algorithms are mainly used for classification problems.

In a regression problem, the aim is to predict a numerical output, whereas in classification the output is categorical.

SVM can also be used for regression problems by using a specific formulation called Support Vector Regression (SVR). However, the emphasis and method used in regression algorithms are different compared to SVM and Random Forest.

Four different datasets were used by three studies [22, 21, 26]. The results span 59.21% to 94% accuracy (F1). Only one study [21] discussed port scanning in detail and included source code for the algorithm.

Table 4: Regression Algorithm Performance

| Authors | Accuracy (F1) | Dataset |
| --- | --- | --- |
| Tojeiro et al. | 94.00 | CICIDS2017 |
| Sirisha et al. | 74.05 | NSLKDD |
| Sirisha et al. | 59.21 | CICIDS2017 |
| Bertoli et al. | 70.00 | MAWILab |
| Bertoli et al. | 92.00 | Bonafide |

### 4.2.4  Naive Bayes

Naive Bayes is a probabilistic algorithm that is good for classification problems, especially when the assumption of independence between features holds. It is fast and simple to implement and can handle large datasets well.

Naive Bayes is different from regression, SVM, and Random Forest algorithms in that it makes a probabilistic prediction based on Bayes' theorem and the assumption of independence between features, whereas the other algorithms make predictions based on a boundary or a combination of trees. In comparison, regression, SVM, and Random Forest algorithms work well for more complex problems where the relationship between features is not necessarily independent and the decision boundary is not clear.

Bakaletz [27] used a custom generated dataset for training and evaluation. The author used two nmap scan techniques- aggressive (NMAP-A) and stealth (NMAP-S). There were five additional datasets used in three [19, 22, 21] out of the remaining four studies [24] in this category. Training and evaluation of NB algorithms demonstrated performances in the range of 55% to 99.7%.

Table 5: Naive Bayes Algorithm Performance

| Authors | Accuracy (F1) | Dataset |
| --- | --- | --- |
| Al-Haija et al. | 99.70 | PSA-2017 |
| Baah et al | 93.02 | CICIDS2017 |
| Bakaletz | 98.00 | NMAP-A |
| Bakaletz | 82.00 | NMAP-S |
| Sirisha et al. | 74.47 | NSLKDD |
| Sirisha et al. | 37.92 | CICIDS2017 |
| Bertoli et al. | 55.00 | MAWILab |
| Bertoli et al. | 78.00 | Bonafide |

### 4.2.5  Neural networks

Neural networks (NNs) are good for a wide range of tasks including classification, speech recognition, and natural

language processing. They are particularly good for complex and non-linear problems, such as recognizing patterns where traditional algorithms struggle.

Neural networks are different from other machine learning algorithms in that they are based on a modeled structure inspired by the human brain. They consist of multiple interconnected nodes, or artificial neurons, that process information. These neurons are organized into layers and the connections between them are associated with weights that are learned during the training process.

In comparison, algorithms such as RF, SVM, NB, and regression algorithms make predictions based on a clear boundary or a combination of trees or linear relationships, whereas NNs are capable of learning complex relationships between the inputs and outputs through their internal structure. NN algorithms have the ability to learn and make predictions based on the examples they are trained on, which makes them highly flexible. However, they can also be more difficult to interpret and train, and require a large amount of data to achieve good performance.

Furthermore, we differentiate between two types of NN: Artificial Neural Networks (ANNs) and Convolutional Neural Networks (CNNs). Both types of deep learning algorithms used for various tasks. The main difference between ANNs and CNNs is the structure and the way they process information. ANNs are fully connected networks, where each neuron in one layer is connected to all neurons in the next layer. This structure makes ANNs computationally expensive and require a large amount of data to achieve meaningful results.

On the other hand, CNNs are specifically designed for classification tasks. They have a unique structure, consisting of convolutional layers, activation layers, pooling layers, and fully connected layers. Convolutional layers are used to extract features, activation layers apply a nonlinear activation function to the output of the convolutional layer, pooling layers reduce the spatial dimensions of the output, and fully connected layers make the final prediction. This unique structure makes CNNs efficient and effective for classification tasks, as it allows them to learn hierarchical representations of the data and identify different objects and features. In comparison, ANNs are not specifically designed for classification and may not achieve the same level of performance as CNNs for these types of tasks.

We did break our own time bounding constraint in this category because Aamir et al. [7] specifically noted neural networks were not being investigated. Hartpence et al. [28] published their work in 2020, therefore we felt obliged to include the work here.

On that note, three of the four studies in this algorithmic category used non-standard datasets. Hartpence et al. [28] provided immense detail in what port scan types exist in the datasets the authors generated as part of their experiments. The general (GEN) dataset contained typical network traffic with port scans intermingled. The

second dataset contained port scans only (TCP). Lv et al. [29] generated a custom dataset by capturing network traffic while executing nmap full connect scans (NMAP-F). Kirtas et al. [30] executed nmap SYN scans (NMAP-Y) while capturing network traffic for the dataset. No study included code snippets or links to code repositories.

Table 6: ANN Algorithm Performance

| Authors | Accuracy (F1) | Dataset |
|---|---|---|
| Hartpence et al. | 99.99 | GEN |
| Hartpence et al. | 99.31 | TCP |
| SaiKiran et al. | 99.11 | CICIDS2017 |
| Kirtas et al. | 87.88 | NMAP-Y |
| Bertoli et al. | 100.00 | Bonafide |

Table 7: CNN Algorithm Performance

| Authors | Accuracy (F1) | Dataset |
|---|---|---|
| Lv et al. | 99.00 | NMAP-F |
| SaiKiran et al. | 63.52 | CICIDS2017 |
| Singh et al. | 99.94 | CICIDS2017 |
| Henry et al. | 98.73 | CICIDS2017 |

## 5   Conclusion

We posed four research questions to guide this systematic review of port scan detection literature. We discovered five algorithms present in the literature: Random Forest, Support Vector Machine, Regression, Naive Bayes, and Neural Network. The literature revealed multiple ways to accurately detect port scanning. Bertoli et al. [21] confirmed RF and ANN algorithms are capable of 100% accuracy against a plethora of port scan types and techniques. Sirisha et al. [22] showed the poorest accuracy, 37.92%, with the authors evaluation of the NB algorithm. ANN seemed to be the strongest type of algorithm across all of the sample papers, followed by RF. There were 11 different datasets used in 34 algorithm experiments with the CICIDS2017 dataset being used in 47% of those experiments. These results originated from 14 research studies published since 2021 and a single study coming from 2020. Unfortunately, we were unable to adequately address the fourth research question (i.e., *What port scanning types and techniques were used for evaluation of those algorithms?*) as only a few studies discussed port scanning in any detail.

Overall, we found a variety of existing work included port scanning as a minor point compared to focuses areas such as general intrusion detection, DDoS, malware, botnets, and so forth. We assume any work demonstrating a capability to detect port scanning mentioned such explicitly and thus would be discoverable in our search. Another assumption underlying this work is research indexing. More specifically, we assume existing work has been indexed and thus was discoverable given our search strings. A final assumption present throughout the literature seems to be the stability of port scanning types and techniques. The dominance of the CICIDS2017 dataset in training and

evaluation supports this point. This assumption will continue to be reasonable as long as significant innovations do not occur in the port scanning research.

It is interesting to note many existing papers experiment with more than one machine learning algorithm. The diversity of results within an algorithm category, across the sample papers in the category, is curious. Perhaps this would not be so unexpected if every paper used a different dataset for training and evaluation. Yet, much of the existing work leverages the CICIDS2017 data. As well, The distinct shift to NN algorithms in 2022 is notable. The prominence of NN over the past year suggests NN and its variations represent a viable research pathway going forward. As a final note, we found the vast majority of studies do not include code snippets or links to GitHub repositories.

Accordingly, replication and reproduction to include specific port scanning techniques with packet captures and algorithm source code would be beneficial. This would be significant because doing so would fill in an existing gap and also enable adjacent research in cybersecurity such as offensive cyber, network security, and so forth. At the same time, such future work might look to incorporate green compute measurements given the increased focus on sustainability in algorithm research. Foundational compute resource metrics can be taken from Bertoli et al. [21] who detailed the compute resource profiles associated with training and evaluating the ML algorithms in their work.

A final area for future exploration is *nature inspired* or *artificial life* (Alife) algorithms. Greensmith et al. [31] showed how a dendritic cell algorithm can be used to detect port scanning. The authors show a theoretical model with pseudocode examples. More recently, Liu et al. [32] extended the dendritic cell concept albeit without evaluation against port scan datasets. Such algorithms should be evaluated using the datasets utilized in the review sample papers and additional Alife algorithms investigated.

## References

[1] Statista Research Department. Most common action varieties in data breaches worldwide in 2019, 2023.

[2] Weijie Wang, Baijian Yang, and Yingjie Victor Chen. Detecting subtle port scans through characteristics based on interactive visualization. In *Proceedings of the 3rd annual conference on Research in information technology*, pages 33–38, 2014.

[3] Habib Ullah Baig and Farrukh Kamran. Detection of port and network scan using time independent feature set. In *2007 IEEE Intelligence and Security Informatics*, pages 180–184. IEEE, 2007.

[4] Tarun Yadav and Arvind Mallari Rao. Technical aspects of cyber kill chain. In *International symposium on security in computing and communication*, pages 438–452. Springer, 2015.

[5] LT Heberlein, GV Dias, KN Levitt, B Mukherjee, J Wood, and D Wolber. A network security monitor. In *Proceedings. 1990 IEEE Computer Society Symposium on Research in Security and Privacy*, pages 296–304. IEEE, 1990.

[6] Azriel Henry, Sunil Gautam, Samrat Khanna, Khaled Rabie, Thokozani Shongwe, Pronaya Bhattacharya, Bhisham Sharma, and Subrata Chowdhury. Composition of hybrid deep learning model and feature optimization for intrusion detection system. *Sensors*, 23(2):890, 2023.

[7] Muhammad Aamir, Syed Sajjad Hussain Rizvi, Manzoor Ahmed Hashmani, Muhammad Zubair, and Jawwad Ahmad. Machine learning classification of port scanning and ddos attacks: A comparative analysis. *Mehran University Research Journal Of Engineering & Technology*, 40(1):215–229, 2021.

[8] Manam Vamsi Krishna and Bhavani Koganti. Machine learning techniques for detecting cyber attacks in networks. *International Journal of Research Sciences and Advanced Engineering*, 9:1–10, 2021.

[9] Gordon Lyon. Remote os detection via tcp/ip stack fingerprinting, Oct 1998.

[10] Marco De Vivo, Eddy Carrasco, Germinal Isern, and Gabriela O De Vivo. A review of port scanning techniques. *ACM SIGCOMM Computer Communication Review*, 29(2):41–48, 1999.

[11] Douglas E Comer and John C Lin. Probing tcp implementations. In *Usenix Summer*, pages 245–255, 1994.

[12] Gary R Wright and W Richard Stevens. *TCP/IP Illustrated, Volume 2 (paperback): The Implementation*. Addison-Wesley Professional, 1995.

[13] Richard J Barnett and Barry Irwin. Towards a taxonomy of network scanning techniques. In *Proceedings of the 2008 annual research conference of the South African Institute of Computer Scientists and Information Technologists on IT research in developing countries: riding the wave of technology*, pages 1–7, 2008.

[14] Elias Bou-Harb, Mourad Debbabi, and Chadi Assi. Cyber scanning: a comprehensive survey. *Ieee communications surveys & tutorials*, 16(3):1496–1519, 2013.

[15] Shanto Roy, Nazia Sharmin, Jaime C Acosta, Christopher Kiekintveld, and Aron Laszka. Survey and taxonomy of adversarial reconnaissance techniques. *ACM Computing Surveys (CSUR)*, 2022.

[16] Harris M Cooper. *Synthesizing research: A guide for literature reviews*, volume 2. Sage, 1998.

[17] Barbara Kitchenham. Procedures for performing systematic reviews. *Keele, UK, Keele University*, 33(2004):1–26, 2004.

[18] Mark Petticrew and Helen Roberts. *Systematic reviews in the social sciences: A practical guide*. John Wiley & Sons, 2008.

[19] Qasem Abu Al-Haija, Eyad Saleh, and Mohammad Alnabhan. Detecting port scan attacks using logistic regression. In *2021 4th International Symposium on Advanced Electrical and Communication Technologies (ISAECT)*, pages 1–5. IEEE, 2021.

[20] Akram QM Algaolahi, Abdullah A Hasan, Amer Sallam, Abdullah M Sharaf, Aseel A Abdu, and Anas A Alqadi. Port-scanning attack detection using supervised machine learning classifiers. In *2021 1st International Conference on Emerging Smart Technologies and Applications (eSmarTA)*, pages 1–5. IEEE, 2021.

[21] Gustavo De Carvalho Bertoli, Lourenço Alves Pereira Júnior, Osamu Saotome, Aldri L Dos Santos, Filipe Alves Neto Verri, Cesar Augusto Cavalheiro Marcondes, Sidnei Barbieri, Moises S Rodrigues, and José M Parente De Oliveira. An end-to-end framework for machine learning-based network intrusion detection system. *IEEE Access*, 9:106790–106805, 2021.

[22] Aswadati Sirisha, Kosaraju Chaitanya, KVSSR Krishna, and Satya Sandeep Kanumalli. Intrusion detection models using supervised and unsupervised algorithms-a comparative estimation. *Journal homepage: http://iieta. org/journals/ijsse*, 11(1):51–58, 2021.

[23] Mahsa Mohseni and Jafar Tanha. A density-based undersampling approach to intrusion detection. In *2021 5th International Conference on Pattern Recognition and Image Analysis (IPRIA)*, pages 1–7. IEEE, 2021.

[24] Emmanuel Kwesi Baah, Steven Yirenkyi, Dominic Asamoah, Stephen Opoku Oppong, Edward Opoku-Mensah, Benjamin Tei Partey, Anthony Kingsley Sackey, Oliver Kornyo, and Evans Obu. Enhancing port scans attack detection using principal component analysis and machine learning algorithms. In *International Conference on Frontiers in Cyber Security*, pages 119–133. Springer, 2022.

[25] N SaiKira, Pradeep Naidu PDS, K Harshini, M Venkateswarlu, et al. Detection of cyber attacks in network using machine learning techniques. *South Asian Journal of Engineering and Technology*, 12(3):138–145, 2022.

[26] Carlos Alexandre Carvalho Tojeiro, Carlos De Jesus Reis, Kelton Augusto Pontara Da Costa, and Thiago José Lucas. Port scan identification through regression applying logistic testing methods to balanced data. *Research Square*, 2022.

[27] Rachel Bakaletz. *A Machine Learning Approach for Reconnaissance Detection to Enhance Network Security*. PhD thesis, East Tennessee State University, 2022.

[28] Bruce Hartpence and Andres Kwasinski. Combating tcp port scan attacks using sequential neural networks. In *2020 International Conference on Computing, Networking and Communications (ICNC)*, pages 256–260. IEEE, 2020.

[29] Conglei Lv, Xiwang Li, and Wei Wang. Application of convolution neural network in network abnormal traffic detection. In *2022 11th International Conference of Information and Communication Technology (ICTech))*, pages 165–169. IEEE, 2022.

[30] M Kirtas, N Passalis, D Kalavrouziotis, D Syrivelis, P Bakopoulos, N Pleros, and A Tefas. Early detection of ddos attacks using photonic neural networks. In *2022 IEEE 14th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP)*, pages 1–5. IEEE, 2022.

[31] Julie Greensmith, Uwe Aickelin, and Steve Cayzer. Detecting danger: The dendritic cell algorithm. *arXiv preprint arXiv:1006.5008*, 2010.

[32] Gang Liu and Jing Wang. Dendrite net: a whitebox module for classification, regression, and system identification. *IEEE Transactions on Cybernetics*, 2021.