

A Comprehensive Survey of Databases and Deep Learning Methods for Cybersecurity and Intrusion Detection Systems

Dilara Gümüşbaş, *Student Member, IEEE*, Tulay Yıldırım , *Member, IEEE*, Angelo Genovese , *Member, IEEE*, and Fabio Scotti , *Senior Member, IEEE*

Abstract—This survey presents a comprehensive overview of machine learning methods for cybersecurity intrusion detection systems, with a specific focus on recent approaches based on deep learning (DL). The review analyzes recent methods with respect to their intrusion detection mechanisms, performance results, and limitations as well as whether they use benchmark databases to ensure a fair evaluation. In addition, a detailed investigation of benchmark datasets for cybersecurity is presented. This article is intended to provide a road map for readers who would like to understand the potential of DL methods for cybersecurity and intrusion detection systems, along with a detailed analysis of the benchmark datasets used in the literature to train DL models.

Index Terms—Cybersecurity, deep learning (DL), intrusion detection system (IDS).

I. INTRODUCTION

CYBERSECURITY systems have been of great importance since the beginning of the computer network era. However, security attacks emerged even before that: in 1941, Alan Turing cracked the Enigma machine, which was designed to cipher messages [1]. Similar incidents have continued to occur through the present day; for example, the electrical system of Massachusetts Institute of Technology (MIT) was hacked in 1950, Yahoo accounts were stolen in 2014, and several worldwide organizations were affected by the WannaCry worm in 2017 [2], [3]. According to the August 2019 threat reports from McAfee Labs, the top ten attack vectors at present are malware, account hijacking, unknown, vulnerability, unauthorized access, targeted attack, code injection, denial of service (DoS), defacement, and theft [4]. Hence, not only known but also unknown attack vectors currently pose a significant cyber threat.

The recent increase in the volume of data generated and transmitted over the Internet, the need to manage the security of such data [5]–[7], and the continuing changes/evolution in intrusion types are causing both the academic and industrial

communities to show increasing interest in the deployment of cybersecurity systems [8], [9]. To shield computer networks from attacks, intrusion detection systems (IDSs) are being deployed to support user authentication, ensure safe access, and prevent loss of privacy. An IDS first collects and processes data and then applies a detection mechanism to raise alarms, which are sent to a human network analyst for further screening.

Different IDSs can employ diverse algorithms for detecting attacks. These algorithms can be classified into the following three categories [10].

- 1) Rule-based algorithms, which use prior knowledge of attacks, such as the corresponding data distributions, to create a rule system and perform detection.
- 2) Statistics-based algorithms, which detect anomalies by building a statistical distribution of intrusion patterns.
- 3) Machine learning (ML) based approaches, in which learning algorithms are adopted to train classifiers that can distinguish among different types of attacks.

Rule-based methods, while simple and fast to execute, cannot compensate for incomplete or noisy data and are difficult to update. To overcome these problems, statistics-based approaches have been proposed to enable the processing of imprecise information; however, such methods entail a high computational cost and have a limited ability to handle large quantities of data. Recently, ML-based approaches have increasingly been studied due to their ability to use complex inference models that can be trained on large quantities of data to detect complex intrusion patterns [11].

Due to the increasing quantities of data transmitted over the Internet, which are leading to the introduction of new networking paradigms (e.g., the Internet of Things (IoT), cloud computing, and fog/edge computing [12], [13]) and complex inference models (e.g., deep learning (DL) [14], [15]), throughout the remainder of this article, we will focus on ML-based approaches to cybersecurity and IDSs.

A. Previous Surveys

This section introduces previous surveys published in the literature on cybersecurity, with a specific focus on ML-based methods. These surveys were chosen based on the criterion of being either the most cited or a pacesetter review on a specific topic. In this article, in contrast to other surveys in the literature,

Manuscript received November 27, 2019; revised March 20, 2020 and April 28, 2020; accepted May 3, 2020. (Corresponding author: Angelo Genovese.)

Dilara Gümüşbaş and Tulay Yıldırım are with the Department of Electronics and Communication Engineering, Yıldız Technical University, Istanbul 34349, Turkey (e-mail: f0415058@std.yildiz.edu.tr; tulay@yildiz.edu.tr).

Angelo Genovese and Fabio Scotti are with the Department of Computer Science, Università degli Studi di Milano, Milan 20133, Italy (e-mail: angelo.genovese@unimi.it; fabio.scotti@unimi.it).

Digital Object Identifier 10.1109/JSYST.2020.2992966

we summarize existing surveys based on their strong points. The objective is to help readers find further material in accordance with their interests.

One of the most cited surveys in the literature is presented in [11]. This survey addresses the different ML methods used in IDSs; describes the structure of Internet Protocol (IP) traffic features, such as port-based, payload-based, and statistical features; and provides insight into feature categories such as packet-level and flow-level features. Although this review was performed using a limited number of papers published between 2004 and 2007 and some of the benchmark datasets described are now out of circulation, it presents valuable information on how to discover and extract novel IP-based intrusion patterns/features from network traffic. The review presented in [16], in addition to describing the various ML-based methods of network intrusion detection, focuses on the characteristics of the types of intrusion. Therefore, this review presents how available statistical features can be used and modified for distributed attack detection and the importance of the threshold used to process these types of features.

In contrast to [11] and [16], the survey presented in [17] focuses on the use of ML and data mining (DM) concepts in IDSs. This review includes a clear explanation of ML and DM algorithms introduced in highly cited papers published before 2016 as well as their usage in IDSs. Notably, this review does not include the newest DL methods, such as convolutional neural networks (CNNs); the newest datasets, such as AWID2018 and CICIDS2017; or practical details such as attack frequency and sample size for the benchmark datasets. Nevertheless, this review does consider fuzzy logic, neural networks, genetic algorithms, and rule-based algorithms.

Similar to the survey presented in [17], the work reported in [18] provides a review of ML methods for IDSs, associating different types of attacks with the features that can be used to detect them. In particular, the associated features can provide insight into how similar features of different types of intrusion can support similar approaches to attack detection. For example, the duration and service features from the KDD99 dataset are the most highly contributing features for detecting both user-to-root (U2R) and remote-to-local (R2L) attacks, often causing these two attack types to be misclassified as one another. Although this article fails to investigate the newest DL algorithms and attack types and their most related features, it provides an extensive survey of feature selection methods.

The review published in [19] surveys ML-based intrusion detection methods alongside newer DL-based methods. Although this survey focuses on certain specific ML and DL methods, such as deep belief networks (DBNs) and recurrent neural networks (RNNs), as well as known benchmark datasets, it does not cover other DL algorithms, such as CNNs, or benchmark datasets such as CICIDS2017. The reviews presented in [14], [20], and [21] also consider DL-based methods. However, they focus on only a subset of these methods, do not discuss benchmark datasets, or do not provide detailed descriptions of the accuracies achieved using DL methods.

In contrast to the previously mentioned surveys, the work presented in [22] focuses on the different types of attacks rather than algorithms for IDSs, without providing details on accuracy.

Furthermore, this article presents an attack taxonomy to provide detailed definitions of various attack types, including how and in which layers they occur. Attack tools are also explained in great detail for readers who wish to build IDSs for protection against specific attack types. Although this article does not provide detailed information about new benchmark datasets or DL algorithms, a brief review on industrial IDSs, such as programmable logic controller systems, is presented. Similarly, the review published in [23] addresses only application-layer distributed DoS (DDoS) attacks, describing how they are hidden behind low traffic and the features used to detect DDoS attacks occurring in the application layer. Furthermore, this review discusses defense mechanisms for protecting against these attacks, such as user puzzles; the limitations of attempts to detect these attacks; and attack generation scenarios.

Finally, there are several surveys that address specific aspects or applications of IDSs. For example, the work reported in [24] focuses on IDSs for IoT systems, describing their taxonomy and placement strategies. In a similar manner, the review presented in [25] discusses DM concepts with IoT applications. Another example is the survey in [26], which covers only unsupervised methods used in IDSs. Although this review is limited to unsupervised methods, it is a good reference for learning about a variety of feature selection methods. Additionally, datasets and EU standards (e.g., the General Data Protection Regulation) for data collection and protection are addressed in this review. Other reviews considering specific aspects of this field include the work described in [27], which focuses on hardware techniques for IDS implementation; the paper presented in [28], which considers only immunity-based approaches; and the survey published in [29], which describes network security techniques for supervisory control and data acquisition systems.

B. Contributions

This work is intended to serve as an extensive survey of databases and methods based on ML and DL that have been introduced thus far in the literature on cybersecurity and intrusion detection. This survey focuses on papers published after 2013, with some exceptions being trendsetter algorithms or highly cited papers.

Compared to the other surveys on intrusion detection discussed in Section I-A, this survey makes the following three main contributions.

- 1) It summarizes previous surveys with regard to their level of detail in describing methods for cybersecurity, with the purpose of encouraging further reading based on the readers' interests.
- 2) It focuses on a practical perspective when describing the relevant datasets, specifically addressing the number of features, the feature types, and attack distributions rather than describing general details, feature selection methods, and algorithms, which are analyzed in other surveys.
- 3) It presents a comprehensive investigation of the newest DL methods for intrusion detection, analyzing their detection capability, performance, and limitations as well as the databases used. This review does not consider previous

TABLE I
LIST OF ACRONYMS AND NOTATIONS USED IN THIS ARTICLE

Notation	Description
ML	Machine Learning
DL	Deep Learning
DM	Data Mining
IDS	Intrusion Detection System
IoT	Internet of Things
IP	Internet Protocol
TCP	Transmission Control Protocol
UDP	User Datagram Protocol
GDPR	General Data Protection Regulation
PCAP	Packet CAPture
SSH	Secure Shell
FTP	File Transfer Protocol
SQL	Structured Query Language
SYN	TCP packet used to request a connection
DoS	Denial of Service
DDoS	Distributed Denial of Service
U2R	User-to-Root
R2L	Remote-to-Local
XSS	Cross-Site Scripting
k-NN	k-Nearest Neighbor
ANN	Artificial Neural Network
SVM	Support Vector Machine
RBM	Restricted Boltzmann Machine
DBN	Deep Belief Network
AE	Autoencoder
CNN	Convolutional Neural Network
RNN	Recurrent Neural Network
LSTM	Long Short-Term Memory
GAN	Generative Adversarial Networks
PCA	Principal Component Analysis

types of ML methods since they have been thoroughly addressed in other survey papers [11], [17], [18].

The remainder of this article is organized as follows. Section II presents a review of cybersecurity datasets, including the data collection steps, feature and attack types, benchmark datasets, and reliability criteria. Section III reviews and analyzes DL-based intrusion detection methods, considering DBNs, autoencoders (AEs), CNNs, long short-term memory (LSTM) networks, and generative adversarial networks (GANs). Section IV provides a discussion of and insights into the limitations and current research trends regarding public datasets and IDSs. Finally, Section V concludes this work. Table I summarizes the acronyms and notations used in this article.

II. CYBERSECURITY DATASETS

This section presents a review of cybersecurity datasets, outlining the data collection steps, feature and attack types, available benchmark databases, and reliability criteria.

A. Data Collection

This section presents the methods of data collection used in cybersecurity applications. Specifically, data collection can be performed in two different ways. The first is based on processing system calls (system logs) from host-based operating systems. The second is based on packet headers and payloads extracted from network traffic packages and from applications using the Transmission Control Protocol (TCP)/IP communication stack [30].

The two main methodologies used to collect network traffic in the second way are full Packet CAPture (PCAP) and the NetFlow protocol.

- 1) *PCAP* enables the collection of the most detailed data from a network because it involves the extraction of

TABLE II
PROGRAMS USED TO CAPTURE AND PREPROCESS NETWORK TRAFFIC

Method	Step	Program	Ref.
PCAP	Capture	libPCAP	[32]
		winPCAP	[33]
		SNORT	[34]
	Preprocessing	Wireshark	[35]
		tshark	[36]
		tcpdump	[37]
		networkminer	[38]
NetFlow	Capture/Preprocessing	rapidminer	[39]
		scapy	[40]
		Cisco NetFlow	[41]
		nfdump	[42]

whole network packets (including packet headers) for all information being transmitted. In particular, the data collected from such packets include the packet size, protocol types, headers of flows, flags, source and destination IP addresses, and source and destination port numbers [31]. However, the information contained in the payload of a packet may be deleted or anonymized due to privacy issues. In fact, a packet payload may contain sensitive data such as private information, instant messaging conversations, or a history of visited websites. In most cases, a tradeoff must be established between anonymizing payloads to protect user privacy and using all collected data to achieve accurate attack detection. This tradeoff is especially important to consider in the case of nonflooding attack types such as R2L and U2R attacks, which are performed using packet payloads.

- 2) *NetFlow* enables the collection of summary information or certain predefined attributes related to the flow of packets in a network. Examples of the features that can be extracted include the number of packets in given a time period or the size of data transmitted over the network. Although data collection via NetFlow is more memory efficient than data collection via PCAP, only summary data are considered, and it is not possible to extract new types of features to address new needs.

The most commonly used programs for performing PCAP are libPCAP, winPCAP, and SNORT. In addition, several programs allow the preprocessing of PCAP files to extract different types of features. For example, such preprocessing programs include Wireshark, tshark, tcpdump, networkminer, rapidminer, and scapy. The most commonly used programs for capturing and preprocessing NetFlow data are Cisco NetFlow and nfdump. Table II summarizes the various programs used to capture and preprocess network traffic using the PCAP and NetFlow methodologies.

B. Feature Types

This section examines the types of features extracted from the available datasets. Although new features are added when novel attack patterns are discovered, there are several reoccurring feature types in the literature.

First, a distinction can be drawn between host-based and network-based data based on the procedure used to collect the data, as described in Section II-A. In most cases, host-based

data are composed of system/operation logs, which consist of attributes such as system calls. Feature extraction from system calls is generally performed using methods based on natural language processing, such as n -grams [43].

On the other hand, network-based data are obtained by collecting network traffic data. However, network traffic is composed of many individual packets/frames, and feature extraction must be performed for each traffic session, known as flow-level traffic data, to reduce the dimensionality of the data and detect intrusions. Such feature extraction is conducted based on three different types of features: basic, traffic-based, and content-based features.

- 1) *Basic features* are extracted from TCP/IP connections and can be classified as header-based, flow-based, connection-based, or packet-based features. Header-based features are related to the packet header and include the source and destination IP addresses, the TCP and User Datagram Protocol (UDP) source and destination ports, the IP protocol, the service, and the IP header length. Flow-based features include attributes computed through analysis of the flow. In particular, a flow is defined as a set of packets having a common set of properties (flow keys), which may include IP addresses, port numbers, or meta-information [44]. Examples of flow-based features are statistical aggregations (e.g., average, maximum, minimum) on the size, time of arrival, and number of inbound/outbound packets in a given time period, the duration of that period, and the type of packets. Connection-based features are related to a particular connection, which is defined as a stream of packets between two specific IP addresses. Such features include the interval between packets, the timestamp, and the time to live. Finally, packet-based features are related to the transmitted data and include the payload and mean number of bytes of a packet. The main advantage of basic features is that they are general and can be used to detect several kinds of attacks [45], [46].
- 2) *Traffic-based features* are associated with either a specific time interval (e.g., 2 s) or a specific number of connections (e.g., 100 connections). These features can be extracted by considering either the same host or the same service. In the first case, the extracted features include statistical sums of connections with the same destination host, whereas in the second case, the extracted features comprise statistical sums of connections to the same service for a fixed amount of time or number of connections [45]. One drawback of traffic-based features is that some attack types span time intervals longer than 2 seconds or a number of connections greater than 100. Examples of such attack types include low-frequency attack types such as U2R, R2L, and low-rate DoS attacks, in which the frequency of the transmitted information is similar to that of legitimate traffic, in contrast to high-frequency attack types, which exhibit a higher frequency than normal traffic. Although some newly proposed connection-based features span time intervals longer than 2 seconds, these features are not fully adequate for identifying such attack patterns [45].

- 3) *Content-based features* are extracted from information embedded in different data portions of packets and include the number of requests, the request type, and the number of failed login attempts. Content-based features are especially useful for detecting low-frequency attack types, which do not exhibit sequential patterns as high-frequency attacks do. In fact, while traffic-based features can be used to detect high-frequency attacks, low-frequency attacks are difficult to detect using only basic and traffic-based features, and in most cases, content-based features are also required [45].

C. Attack Types

This section outlines the various attack types considered in IDSs. In particular, we present the following attack types since they are the ones considered in the most frequently used benchmark datasets [48].

- 1) *DoS* [45] attacks are based on temporarily blocking the normal use of network utilities by flooding the network with traffic. Examples of DoS attacks include botnet, Slowloris, smurf, and SYN flood attacks.
- 2) *DDoS* [46] attacks are based on flooding the server and making it unable to respond by overloading it with service requests. Unlike in DoS attacks, the flooding is performed via many sources. Examples of DDoS attacks include local area network denial (LAND), ping-of-death, RUDY, and teardrop attacks.
- 3) *U2R* [45] attacks involve behaving as a normal user with the aim of detecting system vulnerabilities and gaining root access. Examples of U2R attacks include buffer overflow, rootkit, Perl, and loadmodule attacks.
- 4) *R2L* [45] attacks attempt to use a remote system to gain unauthorized access to and damage the target system. R2L attacks may be combined with U2R attacks, making these types of attacks difficult to differentiate. Examples of R2L attacks include secure shell (SSH) brute force, warezmaster, multihop, imap, and spy attacks.
- 5) *Probe* [45] attacks are based on searching for vulnerabilities throughout the whole network by sending scan packets and gaining information about the system. Examples of probe attacks include Satan, IP sweep, and port sweep attacks.
- 6) *Password* [18] attacks attempt to gain unauthorized access to the system by using guessing techniques to steal passwords. Examples of password attacks include brute force FTP-Patator and brute force SSH-Patator attacks.
- 7) *Injection* [47] attacks use scripts that inject commands/queries with the purpose of gaining unauthorized access and stealing information. Examples of injection attacks include SQL injection and cross-site scripting (XSS).

Table III lists the attack types considered in the most frequently used benchmark datasets, along with their definitions.

Although the definitions provided in Table III can be used to distinguish the different attacks, three additional factors must

TABLE III
ATTACK TYPES REPRESENTED IN THE MOST FREQUENTLY USED CYBERSECURITY BENCHMARK DATASETS

Attack name	Examples	Description
Denial of Service (DoS) [45]	Botnet, Slowloris, smurf, SYN flood	Temporarily blocks the normal use of network utilities by flooding the network with traffic.
Distributed DoS (DDoS) [46]	LAND, ping of death, RUDY, teardrop	Floods the server and makes it nonresponsive to users by overloading it with service requests. Unlike in DoS attacks, the flooding originates from many sources.
User-to-Root (U2R) [45]	Buffer overflow, rootkit, Perl, loadmodule	Behaves as a normal user with the aim of detecting system vulnerabilities and gaining root access.
Remote-to-Local (R2L) [45]	SSH brute force, warezmaster, multihop, imap, spy	Gains local access via a remote system and damages the system. May be combined with U2R attacks, thus making these attacks difficult to differentiate.
Probe [45]	Satan, IP sweep, port sweep	Searches for vulnerabilities throughout the whole network via IP addresses by sending scan packets and gaining information about the system.
Password [18]	Brute force FTP-Patator, brute force SSH-Patator	Gains access to the system after stealing passwords by guessing.
Injection [47]	SQL injection, Cross-Site Scripting (XSS)	Uses a script to inject commands/queries to gain unauthorized access and steal information.

be considered when designing an IDS. First, an attack of one type may be the beginning of another attack of a different type. In this case, the characteristics of the true attack will be a combination of the characteristics of both attacks. Second, some attack characteristics may evolve over time. For instance, DDoS attacks are mostly understood to be high-frequency attacks that flood the bandwidth of a network; however, DDoS attacks in the application layer are low-frequency attacks that flood the server instead of flooding the network. Third, some attack types may show similar patterns. For example, both DoS and probe attacks, in most cases, exhibit sequential patterns and involve a large number of connections to the same host, whereas R2L and U2R attacks are both embedded in packets. Therefore, although DoS and probe attacks are easy to differentiate from R2L and U2R attacks, it may not be as easy to differentiate DoS attacks from probe attacks or U2R attacks from R2L attacks due to their similar embedding patterns.

To increase the effectiveness of differentiating among attack types, several studies have investigated which types of features are effective for detecting particular attack types. For example, Mishra *et al.* [18] report that on the basis of the features contained in the KDDCUP99 dataset, even though DoS attacks can be differentiated using basic and traffic-based features, considering some sparse features, such as flags, destination IP addresses, percentages of connections to the same service, and percentages of connections to the same port, can result in more effective detection. Similarly, duration, service, destination host same service rate, and flag features are vital for detecting probe (scanning) attacks. The most important features for detecting U2R attacks are the number of failed logins, number of shells, number of roots, duration, and service. For R2L attacks, the most important features are the duration, service, service bytes, destination bytes, number of failed logins, count, destination host count, and destination host service count. As seen above, the features used to detect attacks of the probe, U2R, and R2L types show a high degree of similarity, which explains why these three attack types are often misclassified among each other.

D. Benchmark Datasets

This section introduces and analyzes benchmark datasets for intrusion detection, considering both the extent to which they reflect novel attack types due to the evolving nature of

intrusion patterns over time and their shortcomings. For the benchmark datasets considered in this section, Table IV lists the most frequently used datasets, whereas Table V summarizes the distribution of the samples in each dataset across the different attack types considered.

1) *AWID2018*: Also known as CSE-CIC-IDS2018, this dataset includes databases for training and testing collected using two different capture procedures. The data collected using the first procedure consist of full-packet network traffic with system logs, whereas the data collected using the second procedure consist of reduced packet traffic. The dataset includes two different labels for attacks: a main attack label and a subattack label. This dataset has the advantages of including the newest attack types, such as password attacks based on the SSH/FTP brute force approach, injection attacks based on SQL injection, and flooding attacks based on DoS. However, the data exhibit some limitations, such as noisy, misleading features and uncategorized samples. The dataset consists of 155 features extracted using Wireshark [49].

2) *CICIDS2017*: This dataset was created from realistic traffic data at the Canadian Institute for Cybersecurity of the University of New Brunswick (UNB) in 2017 and includes a full-packet dataset with 152 features and raw PCAP files [50]. The dataset considers attacks and subattacks such as injection attacks based on SQL injection and XSS, password attacks based on brute force FTP-Patator and brute force SSH-Patator, and flooding attacks based on DoS, Goldeneye DDoS, HULK DDoS, slow HTTP DDoS, Slowloris DDoS, and Heartbleed. Although the criteria for a reliable dataset proposed by [54] are satisfied, one feature among the attributes is duplicated.

3) *KDD99/KDDCup99*: Also known as KDDCup99, the KDD99 dataset was created using DARPA 1998 PCAP files and includes full-packet data, divided into subsets for training and testing [51].

This dataset considers DoS-based subattacks such as back, LAND, ping of death, teardrop, Neptune, and smurf attacks; U2R subattacks such as buffer overflow, loadmodule, Perl, and rootkit attacks; R2L subattacks such as ftp-write, t guess-password, imap, multihop, PHF, spy, warezclient, and warezmaster attacks; and probe-based subattacks such as port sweep, IP sweep, NMAP, and Satan attacks. As of 2019, this dataset remains the most widely used benchmark dataset in the field of network intrusion detection. However, this dataset suffers from several limitations, including duplicated samples,

TABLE IV
OVERVIEW OF THE MOST FREQUENTLY USED CYBERSECURITY BENCHMARK DATASETS

Ref.	Name	Year	Num. of features	Num. of samples	Attack types	Separate train-test sets
[49]	AWID2018	2018	155*	210900113 (full) 2326218 (reduced)	Flooding, impersonation, injection	Yes
[50]	CICIDS2017	2017	152***	2830743	DoS/DDoS, port scan, FTP-Patator, SSH-Patator, bot, web attacks, infiltration, Heartbleed	No
[51]	KDD99	1999	42*	4900000 (full) 494021 (subset) 311029 (testing)	DoS, probe, U2R, R2L	Yes
[45]	NSL-KDD	2009	42*	125973 (training) 22544 (testing) 25192 (training) 11850 (testing)	DoS, probe, U2R, R2L	Yes
[52]	Kyoto	2006-2015	24*	Various	Known, unknown	No
[53]	UNSW-NB15	2015	49**	2540047 (full) 175341 (training) 82332 (testing)	Fuzzers, worms, shellcode, analysis, backdoors, DoS, exploits, generic, reconnaissance	Yes

Notes: * = including 1 feature as a label; ** = including 2 features as labels; *** = at the time of this survey.

TABLE V
DISTRIBUTIONS OF ATTACK TYPES IN THE MOST FREQUENTLY USED BENCHMARK DATASETS

Name	AWID2018									
Attack	Normal		Flooding		Impersonation		Injection			
N. samples	205074514		1409392		2361892		2054315			
(Perc.)	(97.24%)		(0.67%)		(1.12%)		(0.97%)			
Name	CICIDS2017									
Attack	Benign	DoS	DDoS	Port scan	FTP-P.	SSH-P.	Bot	Web att.	Infiltr.	Heartb.
N. samples	2273097	252661	128027	158930	7938	5897	1966	2180	36	11
(Perc.)	(80.3004%)	(8.9257%)	(4.5228%)	(5.6144%)	(0.2804%)	(0.2083%)	(0.0695%)	(0.077%)	(0.0012%)	(0.0003%)
Name	KDD99									
Attack	Normal		DoS		Probe		U2R		R2L	
N. samples	972781		3683370		41102		52		1126	
(Perc.)	(20.71%)		(78.4%)		(0.8897%)		(0.0001%)		(0.0002%)	
Name	NSL-KDD									
Attack	Normal		DoS		Probe		U2R		R2L	
N. samples	77054		53385		14077		252		3649	
(Perc.)	(51.9%)		(35.9%)		(9.5%)		(0.2%)		(2.5%)	
Name	Kyoto									
Attack	Normal				Known attacks			Unknown attacks		
N. samples	1186780				11218206			563		
(Perc.)	(9.5706%)				(90.429%)			(0.0004%)		
Name	UNSW-NB15									
Attack	Fuzzers	Worms	Shellcode	Analysis	Backdoors	DoS	Exploits	Generic	Rec.	
N. samples	24246	174	1511	2677	2329	16353	44525	215481	13987	
(Perc.)	(7.572%)	(0.054%)	(0.47%)	(0.83%)	(0.72%)	(5.112%)	(13.8%)	(67.092%)	(4.35%)	

Notes: N. samples = Number of samples; Perc. = Percentage; FTP-P. = FTP-Patator; SSH-P. = SSH-Patator; Web att. = Web attacks; Infiltr. = Infiltration; Heartb. = Heartbleed; Rec. = Reconnaissance. The largest remainder method was used when computing the percentages to ensure a total of 100%.

different probability distributions between the training and test data, unbalanced classes, and a lack of coverage of the newest attack types.

4) *NSL-KDD*: This dataset was created by erasing all duplicate records from the KDD99 dataset and using sampling techniques to balance the number of data samples in each class [45]. This dataset includes separate databases for training and testing, where the test database consists of 14 subattack types that are not present in the training database. NSL-KDD is not subject to most of the limitations of the KDD99 dataset; however, this dataset still lacks newer attack types.

5) *Kyoto*: This dataset was created from honeypots at Kyoto University and consists of traffic data collected daily between 2006 and 2015 [52]. The dataset includes 24 features, 14 of which are in common with the KDD99 dataset, and labels indicating normal data, known attacks, and unknown attacks. The dataset is missing data from some days and months during the time of its collection, and the average number of samples per month is approximately 12 million. Since the traffic was captured from honeypots, which are designed to protect against less advanced attackers, most of the monitored attacks did not originate from advanced attackers. Therefore, the dataset may not be representative of realistic attacks.

6) *UNSW-NB15*: This dataset was synthetically created at the Cyber Range Lab of the Australian Centre for Cybersecurity and includes full, training, and test datasets as well as raw PCAP files. The dataset includes 49 features and two label attributes: the first label describes the attack, and the second label is binary. The dataset considers attacks such as fuzzers, backdoors, shellcode, DoS attacks, worms, generic attacks, reconnaissance attacks, exploits, and analysis attacks [53]. One of the limitations of this dataset is the existence of several missing samples.

7) *DARPA*: This dataset was created at the MIT Lincoln Laboratory in 1998 and includes full, training, and test sets of raw PCAP files [55]. The newer versions of the DARPA dataset, DARPA 1999 and DARPA 2000, are based on the 1998 version. This dataset is one of the most commonly used intrusion detection datasets; however, it is commonly considered to be outdated and to contain irregularities [56].

8) *ISCX IDS 2012*: Also known as UNB or UNB ISCX 2012, this dataset was created at UNB in 2012 and includes full-packet network data [57]. The dataset includes normal traffic data and attack data for attack types such as infiltration, DoS, DDoS, and brute force SSH attacks. Although this dataset includes some of the newest attack types, it is criticized as being unrealistic for not containing sufficient Internet background noise, as it

consists of pure network traffic rather than data received by any real device [58].

9) *CIC DoS*: This dataset was created at the Canadian Institute for Cybersecurity of UNB in 2017 [59]. It considers the application layer and incorporates data that describe high-volume (traditional) DoS attacks, data corresponding to low-volume DoS attacks, and normal data from the ISCX IDS 2012 dataset.

10) *Gure-KddCup*: This dataset was created using the PCAP data from the DARPA 1998 dataset [60]. It includes features similar to those of the KDD99 dataset, with the addition of payload information and other new features, such as IP addresses and port numbers, to make U2R and R2L attacks more visible/distinguishable [61].

11) *Cyber Defence Exercises (CDX)*: The CDX dataset [62] was collected from the United States Military Academy network in 2009 and consists of PCAP data extracted from system logs, divided into intrusion traffic and normal traffic [56].

12) *ASN-CDX*: This dataset was created from the CDX network traffic data in 2009. The dataset includes 5772 samples, each with 875+1+1 features. It includes distributed features often used in detecting low-frequency attacks, such as the number of packets and the total bytes in/out from 4 to 54 s. In some cases, the features have been converted with the fast Fourier transform to increase their discriminative ability. This dataset has two attack label attributes: the first label discriminates between legitimate and malicious traffic, and the second label indicates whether the attack is based on buffer overflow. However, this dataset lacks traffic diversity since it consists only of buffer overflow attacks [63].

13) *Lawrence Berkeley National Laboratory (LBNL)*: This dataset was created at the LBNL between 2004 and 2005. Although the dataset includes packet headers, the payloads are anonymized due to privacy issues, which limits its informativeness [64].

14) *ISOT*: This dataset was created in 2010 by combining Storm, Waledac, and Zeus botnet attack data from the French Chapter of the HoneyNet Project and normal traffic data from the Traffic Lab at Ericsson Research and LBNL [65].

15) *MAWI*: This dataset was collected by the MAWI Working Group in Japan and includes continuously updated traffic data from 2001 to 2019. A graph-based methodology has been used to label the raw data as either abnormal or normal [66]. One of the limitations of this dataset is duplicated packets.

16) *CTU-13*: This dataset is a combination of botnet traffic data, normal data, and background data collected at Czech Technical University in Prague (CTU) in 2011. Although the data consist of a variety of botnet scenarios and extended truncated versions of PCAP files with complete TCP, UDP, and Internet Control Message Protocol headers, the dataset is specifically designed only for botnet detection. Therefore, it is considered unrealistic to mix these data with normal and background traffic [67].

17) *UMass*: This dataset was collected between 2004 and 2018 and contains traffic data such as Tor traffic data, Gateway Link 3 Trace data, web requests, and response data. However, most of the data were collected under similar network traffic conditions and lack a broad variety of attacks [68].

18) *Twente*: This dataset was created from honeypots at the University of Twente in 2009 and consists of more than 14 million flows and more than seven million alerts. In this dataset, some samples are left unlabeled, and informative data from the packet headers and payloads are anonymized [69]. This dataset has the limitation that traffic originating from honeypots does not represent realistic attacks since honeypots are designed to protect against less advanced attackers.

19) *CAIDA*: The CAIDA dataset consists of a variety of different databases that are specific to particular events, such as network telescope and DDoS databases [58], [70]. Although there are a few up-to-date databases, such as CAIDA DDoS, most do not accurately represent the different possible types of attacks. For instance, the DoS attack databases consist only of spoofed-source DoS attacks and exclude other versions of DoS attacks.

20) *DEFCON*: The DEFCON datasets are created for intrusion modeling competitions held every year. Although these datasets are continuously created, they focus only on intrusions and attacks and lack normal background traffic [58]. Therefore, they are not frequently used for network intrusion detection.

21) *Others*: In addition to the most commonly used benchmark datasets, a variety of publicly available raw traffic datasets exist. These datasets include Metrose, UNIBS 2009, TUIDS, the University of Napoli traffic dataset, payload datasets such as the CSIC 2010 HTTP Dataset, the UNM system call dataset, and an enormous variety of network traffic from the Capture the Flag Competitions and CDX. Moreover, several host-based datasets also exist, including the ADFA Linux Dataset (ADFA-LD), the ADFA Windows Dataset (ADFA-WD), and the ADFA Windows Dataset Stealth Attacks Addendum (ADFA-WD:SAA) [71].

III. DL-BASED INTRUSION DETECTION METHODS

Traditional ML-based methods for cybersecurity include approaches based on the k -nearest neighbor (k -NN) algorithm, k -means clustering, artificial neural networks (ANNs), fuzzy logic, Bayesian networks, hidden Markov models, self-organizing maps, decision trees, evolutionary classifiers, support vector machines (SVMs), and rule-based systems [17], [18], [22], [26]. In this survey, we focus on the more recent DL-based approaches, which have not been covered in detail in previous surveys.

To provide up-to-date descriptions of the recent methods developed for cybersecurity, this section describes DL-based methods for intrusion detection. For each algorithm, we consider evaluation criteria such as a fast run/convergence time, a high detection ability with a low false positive rate, adaptability to novel intrusions, computational efficiency, and scalability [16]. In the remainder of this section, we consider DL methods based on DBNs, AEs, CNNs, LSTM networks, and GANs [15]. A summary of the presented DL methods in the IDS context is presented in Table VI.

A. Deep Belief Networks

DBNs are a type of ANN obtained by stacking together several restricted Boltzmann machines (RBMs [77]), which act as the layers of the DBN, and introducing connections between the

TABLE VI
SUMMARY OF DL-BASED METHODS FOR INTRUSION DETECTION

Method	Description	Pros	Cons
Deep Belief Networks (DBNs) [72]	Stacks of Restricted Boltzmann Machines (RBMs) with connections between the layers but not within each layer.	Fast and unsupervised layer-by-layer learning in a greedy fashion. Unsupervised dimensionality reduction.	Training uses an approximation of the gradient.
Autoencoders (AEs) [73]	Encoder-decoder structure that maps input data to a hidden space and then reconstructs them.	Can be trained in an end-to-end manner using learning algorithms based on gradient descent. Unsupervised dimensionality reduction.	Requires an additional ML model to perform classification.
Convolutional Neural Networks (CNNs) [74]	Sequences of convolutional layers trained via gradient descent.	Performs classification while automatically learning data representations. Learns discriminant spatial patterns invariant to translation and shifting.	Computationally expensive to train. Not naturally suited to processing data in time-series form.
Long Short-Term Memory (LSTM) [75]	Neurons arranged in a temporal sequence, able to maintain memory for arbitrary intervals of time.	Can natively process time-series data.	The research community is increasingly focusing on CNNs rather than LSTM networks.
Generative Adv. Networks (GANs) [76]	Combination of a generator, which generates data starting from a random distribution, and a discriminator, which distinguishes real data from synthetic data.	Learns data distributions in an unsupervised manner.	Often requires visual inspection of the results.

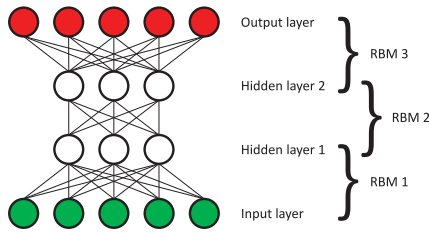


Fig. 1. Example of a DBN. DBNs are obtained by stacking together several RBMs, which act as the layers of the DBN. In DBNs and RBMs, neurons of different layers are fully connected, whereas connections within the same layer are restricted.

layers but not within each layer. The RBMs used to construct a DBN consist of two main layers, one visible and one hidden, constituted by a variable number of neurons. Additionally, within each RBM, the neurons of different layers are fully connected, whereas the connections within the same layer are restricted [72]. Fig. 1 shows an example of a DBN.

Because of their layered structure, DBNs have the advantage that fast learning procedures can be used, which can be applied in a greedy fashion, layer by layer, in an unsupervised way [78]. As a consequence of this advantage, methods based on DBNs were among the first DL-based approaches studied for intrusion detection. In addition, the ability to train DBNs using fast and unsupervised learning algorithms makes them particularly suitable for performing a preliminary dimensionality reduction step, with the aim of extracting a compact and discriminant representation of the data, without the need for labels, even in the case of large intrusion detection databases. For example, in the method described in [79], proposed in 2011, a DBN was first applied to perform a feature reduction step, and an SVM was then used to classify the intrusions contained in the NSL-KDD dataset. Similarly, the approach proposed in [80] is based on a DBN, the parameters of which are first optimized via particle swarm optimization to map the input data to a space of reduced dimensionality. Then, a probabilistic neural network is trained to perform classification.

More recent methods have introduced ML architectures based on deeper DBNs, such as the approaches described in [81]–[83], which have achieved improved detection accuracy on the NSL-KDD and KDD99 datasets.

DBNs have the drawback that it is computationally unfeasible to train them end to end in a supervised way using gradient

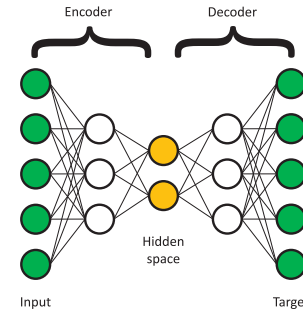


Fig. 2. Example of an AE. The same data are used as both input and target: the encoder maps the input data to a hidden space in a nonlinear manner, and the decoder reconstructs the input data by mapping the encoded data back to the original input space.

descent methods. Due to this drawback, in most cases, DBNs are trained using training algorithms based on contrastive divergence, which rely on an approximation of the gradient [84]. Recently, however, the growing availability of computing power and GPU-based training architectures (e.g., CUDA [85]) has made it possible to train DL models using end-to-end learning algorithms based on gradient descent, without the need to approximate the gradient [15]. Examples of recent DL models trained end to end include AEs, CNNs, LSTMs, and GANs.

B. Autoencoders

AEs are a type of ANN used to learn and reconstruct a representation of input data. In AEs, the same data are used as both input and target, with the purpose of learning a model that can extract a compact and discriminant representation of the input data in an unsupervised manner. Such a representation can then be used as input to a classifier to perform detection. An AE consists of two components: an encoder and a decoder. The encoder maps the input data to a hidden space in a nonlinear manner, and the decoder reconstructs the input data by mapping the encoded data back to the original input space. The purpose of the decoder is to minimize the reconstruction error, defined as the difference between the input data and the reconstructed data [73]. Fig. 2 shows an example of an AE.

Because of their ability to extract a compact and highly discriminative representation with reduced dimensionality from input data, AEs are often used as a preprocessing step in intrusion detection. In most cases, the related approaches presented in

the literature involve using AEs to preprocess the input data, followed by the application of an ML classifier. In contrast to DBNs, AEs can be trained in an end-to-end manner using learning algorithms based on gradient descent, without the need to approximate the gradient [15]. For example, in the method described in [86], an AE with seven layers is used to obtain a compact and discriminant representation of the input data. On the NSL-KDD dataset, this method achieves superior detection accuracy compared with other dimensionality reduction methods based on principal component analysis (PCA) and kernel PCA. Its main limitation is the lack of information about the shallow classifiers used to classify the data after the dimensionality reduction phase. Similar to that in [86], the method proposed in [87] involves applying an AE to the input data to extract features, which are then classified using shallow classifiers such as naive Bayes, k -NN, and SVM classifiers. This work considers the NSL-KDD dataset and reports higher accuracy than that achieved by previous methods, particularly when using the naive Bayes classifier.

Several methods in the literature combine an AE with a density estimation model to achieve greater detection accuracy. For example, the method proposed in [88] adopts a combined approach based on an AE and density estimation. This method achieves a high detection accuracy on the NSL-KDD dataset, especially for DoS and probe attacks. The method described in [89] extends the previous method by combining an AE with a Gaussian mixture model to perform intrusion detection. The model consists of an estimation network, which evaluates the densities of the samples in a low-dimensional space, and a compression network, which projects the data into a lower dimensional space. A procedure based on joint parameter optimization is used to update the model parameters. On the KDD99 dataset, this method achieves a significant accuracy improvement compared to baseline methods using pretrained AEs.

A variant of AEs is represented by sparse AEs, which use a sparsity constraint to further reduce the dimensionality of the obtained representation [73]. Specifically, the method described in [90] uses a sparse AE combined with a softmax regression classifier to perform intrusion detection. The method achieves higher accuracy than previous models on the NSL-KDD dataset but considers only binary classification, differentiating between normal and anomalous traffic.

As the available computational power has increased, recent methods based on AEs have also considered stacked AEs (SAEs) [91], which consist of several AEs trained separately and then “stacked” to obtain a deeper model and a more discriminant representation. The method proposed in [92] uses a model based on an SAE to preprocess raw traffic data in the CTU-13 dataset. Similarly, the method proposed in [93] uses an SAE to process traffic data captured from home wireless networks, representing several types of DDos attacks. To improve the detection accuracy, the method introduced in [94] combines an SAE with a random forest classifier. The method has been tested on the KDD99 and NSL-KDD datasets by reducing the feature dimensionality of the input data and performing five-class classification. The method achieves high overall detection accuracy but exhibits low accuracy in detecting U2R and R2L attacks. Similarly, the

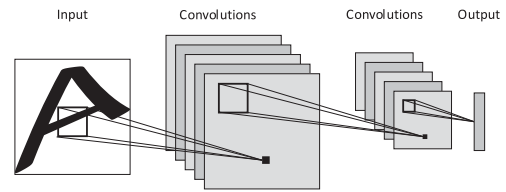


Fig. 3. Example of a CNN. The layers perform subsequent processing by convolving the data with banks of filters.

method described in [95] achieves high detection accuracy on the KDD99 dataset based on a combination of four AEs. Instead of using a random forest classifier, the method presented in [96] combines SAEs with a radial basis function classifier to achieve high detection accuracy on the AWID2018 dataset.

Recently, one of the most commonly used types of AE models has been variational AEs (VAEs) [97]. Their main novelty is that, whereas AEs use a deterministic discriminative model, VAEs use a probabilistic generative model to reconstruct the input data. As a consequence, VAEs are less prone to overfitting than AEs are. In the field of intrusion detection, the method described in [98] combines VAEs with a gradient-based fingerprinting detection model. In this method, gradient-based fingerprints are first extracted from NetFlow data taken from the UGR16 dataset [99], and VAEs are then applied for feature reduction. The work reported in [100] presents a more comprehensive evaluation conducted by combining VAEs with several classifiers, such as naive Bayes, SVM, decision tree, and random forest classifiers. This method achieves good results on the NSL-KDD and UNSW-NB15 datasets, especially when using the decision tree and random forest classifiers.

Despite the ability of AE-based methods to automatically obtain a compact and discriminant feature representation that can be adapted to the input data, AEs have the drawback of requiring an additional ML model to perform classification based on the obtained feature representation. To compensate for this drawback, the use of CNNs is increasingly being considered in recent DL-based methods developed for IDSSs. This is because CNNs can be trained to process input data to automatically learn a compact and discriminant representation [15] while simultaneously classifying the obtained representation into the corresponding attack types [21].

C. Convolutional Neural Networks

CNNs are a type of ANN in which the layers are structured to process input data in the form of multidimensional signals such as images or three-dimensional (3-D) volumes. The different layers in a CNN perform subsequent processing by convolving the data with banks of filters, with parameters typically learned via gradient descent. CNNs have the main advantage of performing classification while automatically learning data representations, without the need for a handcrafted feature extraction step [74], [101]. Due to this advantage, CNNs have been successfully applied in several scenarios [15], [102]–[105]. Fig. 3 shows an example of a CNN

In addition to their advantage of automatically learning data representations, CNNs have been especially successful in processing data with the aim of learning discriminant spatial patterns among the features that are invariant to translation and shifting [101]. In the context of IDSs, the advantages of CNNs have been useful for classifying attack types by considering the relationships among the features while requiring minimal preprocessing of the data [21].

Recently, however, the vast majority of CNN architectures have been structured to process data in the form of images [15]. Therefore, to use a CNN for intrusion detection, a preprocessing step must be performed to transform the features into a 2-D format that can be processed by the CNN. Several methods have been proposed for transforming features into a 2-D format. For example, the preprocessing method proposed in [106] converts feature attributes into binary vectors. The method converts symbolic attributes, such as flag, service, and protocol type attributes, into binary vectors using one-hot encoding [107]. Then, continuous attributes are converted by performing min-max normalization, discretizing the normalized values into ten intervals, and applying one-hot encoding. Finally, the obtained vectors are combined and reshaped to form a 2-D image. Similarly, the preprocessing approach presented in [108] converts malware binaries into grayscale images. A malware file is first read as a vector of 8-b binary numbers, and each binary number is then converted into its equivalent decimal value. Finally, the resulting decimal vector is reshaped into a 2-D grayscale image.

Recent preprocessing methods for CNN-based intrusion detection have extended grayscale image representations to consider multiple channels. For example, the new encoding method proposed in [109] is designed to give equal weight to each feature, producing a feature representation with 24 b for each pixel, similar to an RGB color image.

To accelerate the transformation process, some methods involve performing feature selection before converting the data into an image-based format. For example, in the method described in [110], feature selection is applied via a genetic algorithm.

After the features are transformed into an image-based format, a CNN-based approach is applied to classify the obtained images and perform intrusion detection. Several such approaches have been proposed in the literature. Most of these CNN-based methods rely on a layer structure based on an existing architecture. In particular, several methods use architectures based on LeNet [74], ResNet [111], GoogLeNet [112], or VGG-16 [113]. The LeNet architecture is used in the methods described in [114]–[118], which achieve high detection accuracy, especially on the AWID2018 dataset [119]. Similarly, the ResNet and GoogLeNet architectures are used in the method proposed in [106], which has been tested on the NSL-KDD dataset after preprocessing. Although satisfactory results are achieved on this dataset, detection rates are not reported for each class. The VGG-16 architecture is used in the method presented in [108], which has been applied to the Malign Dataset and the Microsoft Malware Dataset. This method achieves almost the same accuracy as the winner of the Microsoft Malware Dataset Challenge [120].

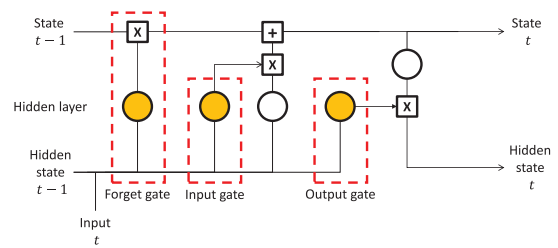


Fig. 4. Example of an LSTM network. The neurons are connected following a temporal sequence. The forget, input, and output gates control which information is preserved in the network and passed to the next time step.

In addition to methods based on existing networks, there are some CNN-based methods for which innovative ML architectures have been proposed. For example, Nguyen *et al.* [121] propose a novel CNN and present its application to the KDD99 dataset. The method achieves a high classification accuracy for five types of DoS attacks, exhibiting performance superior to that of naive Bayes and k -NN classifiers. An innovative architecture is also proposed in [122] based on a convolutional AE; however, this method has been tested only on a private dataset.

Methods based on CNNs have recently achieved high accuracy on several intrusion detection datasets due to their ability to simultaneously learn a compact representation of the input data and perform adaptive classification. However, CNNs are primarily useful for learning discriminant spatial patterns from input data, whereas they are not naturally suited for processing data in the form of time series with the intent of learning discriminant temporal patterns. To overcome this disadvantage, several recent methods have considered the use of LSTM networks, which are specifically structured for learning temporal patterns by processing new data while maintaining a memory of previous samples [75].

D. LSTM Networks

An LSTM network is a type of ANN based on an RNN in which the neurons are connected following a temporal sequence. However, in contrast to traditional RNNs, LSTM networks have a deeper structure of hidden neurons with the ability to maintain a memory of previous inputs for arbitrary intervals of time [75]. Due to this node arrangement, RNNs and LSTM networks are often used to process data in the form of time series [123]. Fig. 4 shows an example of an LSTM network.

The ability of LSTM networks to process time-series data has proven useful in the IDS context since datasets for cybersecurity and intrusion detection are often structured as sequences of features evolving over time. Due to this advantage, several intrusion detection methods in the literature are based on LSTM networks [124]. Among these methods, the approach proposed in [125] applies a three-layer LSTM network. It achieves high detection accuracies on the KDD99, ADFA-LD, and UNM datasets. Similarly, the method proposed in [126] uses a cascade of three LSTM network modules, combined using a voting mechanism, to achieve an increased intrusion detection accuracy.

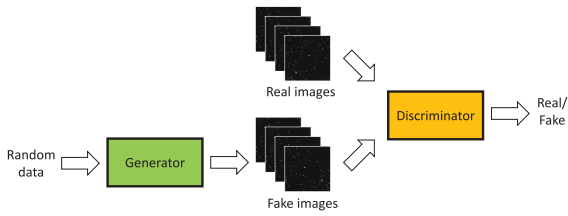


Fig. 5. Example of a GAN, composed of a generator, which generates data starting from a random distribution, and a discriminator, which distinguishes real data from synthetic data. The generator and discriminator are trained in an alternating fashion and, in most recent architectures, have CNN-based structures.

To exploit both the accuracy of LSTM networks in processing time series and the capability of CNNs to extract spatial patterns from images, recent methods have increasingly considered combinations of LSTM and CNN architectures for intrusion detection. For example, the method proposed in [127] uses an LSTM network combined with a CNN to perform multiclass detection of anomalies in the KDD99 dataset. Similarly, the approach described in [128] uses both a CNN and a hybrid LSTM-CNN model to perform detection. In some cases, hybrid LSTM-CNN models have been developed based on existing architectures, such as the method developed in [129], which relies on a CNN designed based on the LeNet model. This method has been tested on recent databases, including CICIDS2017 and CTU-13.

Despite the ability of LSTM networks to natively process time-series data, the introduction of novel and advanced CNN architectures is shifting the attention of the research community toward the use of CNNs in a wider range of application scenarios, including the learning of temporal patterns [130]. In fact, current research trends are increasingly focusing on CNN architectures that are deeper (e.g., ResNet) [111] or lighter in weight (e.g., MobileNet [131]) and on computing platforms specifically designed to accelerate the training of such architectures [85]. Consequently, CNN-based methods tend to outperform models based on recurrent architectures, such as LSTM-based models, in most cases [132].

E. Generative Adversarial Networks

GANs are DL models that can learn and mimic the distribution of input data to generate synthetic samples with a strong resemblance to the original data. Specifically, a GAN is structured as a combination of a generator, which generates data starting from a random distribution, and a discriminator, which distinguishes real data from synthetic data. The generator and discriminator are trained in an alternating fashion until equilibrium is reached [76], [133]. In most recent applications, the generator and discriminator of the GAN are structured as CNNs, with the consequence that recent GANs can generate synthetic image samples with a high degree of realism [134]. Fig. 5 shows an example of a GAN.

GANs have the main advantage of being able to learn the distribution of the input data in an unsupervised manner, that is, without requiring class labels. In the IDS context, this characteristic is useful for learning the characteristics of data distributions

in specific situations (e.g., under normal conditions). Due to this advantage, many recent methods developed for IDSs use GANs trained on existing datasets to detect anomalies, where the training data include only data captured in specific situations. Among these methods, a CNN-based GAN is introduced in [135] to learn the characteristics of data captured under normal conditions. Then, the method is used to detect anomalies by computing the distance between freshly captured data and normal data. To achieve faster detection, an algorithm is proposed in [136] that improves the computational efficiency of the GAN described in [135] while achieving a similar accuracy on the KDD99 dataset.

Despite the ability of GANs to simulate input data distributions, the synthetic data generated from the learned distribution may be insufficiently realistic compared with the real data, and thus, manual (e.g., visual) inspection may be required to achieve good results. In the case of IDSs, such visual examinations of the feature vectors could be relatively difficult to perform compared with cases in which a GAN is used to generate images of known objects or people.

IV. DISCUSSION

This section presents a discussion of the limitations, challenges, and research trends of the current databases and intrusion detection approaches for cybersecurity applications. In particular, we will focus on the issue of dataset reliability and on research directions regarding novel features for intrusion detection.

A. Dataset Reliability

The recent rise in the number of ML-based approaches, particularly those based on DL, has resulted in an increase in the accuracy of intrusion detection that can be achieved using state-of-the-art methodologies. However, the performance of DL-based methods strongly depends on the quantity and quality of the data available [15], with the consequence that the biases and limitations of the datasets used to train the models directly affect the reliability of the predictions.

Currently, although the majority of works focus on researching algorithms that can yield improved detection results, a few studies have been dedicated to evaluate the reliability of benchmark datasets. As a first step toward evaluating dataset reliability, the work proposed in [137] discusses the criteria for a reliable benchmark dataset, which concern the diversity of the traffic data, the diversity of the protocols, the volume of collected data, the diversity of the attacks considered, the inclusion of novel attack types, the inclusion of full payloads without anonymization, the presence or absence of informative features, the updatability, the consideration of realistic traffic, the extent of labeling, and the size of the feature set. Finally, any discussion of dataset reliability should consider the ability of a dataset to adapt to changes over time, for example, by mimicking statistically normal traffic in accordance with upcoming needs.

Similarly, the work described in [54] proposes the following 11 criteria for assessing the reliability of a dataset for intrusion detection:

- 1) attack diversity;
- 2) anonymity;
- 3) available protocols;
- 4) complete capture (with payloads);
- 5) complete interaction;
- 6) complete network configuration;
- 7) complete traffic;
- 8) feature set;
- 9) heterogeneity (all network traffic and system logs);
- 10) correct labeling;
- 11) metadata (full documentation of data collection).

In addition to complying with these criteria, a reliable dataset should also provide a means of anonymizing the payload information to guarantee users' privacy [9].

Among the criteria considered in [54] and [137], attack and traffic diversity play a major role, since a limited diversity or a high imbalance among attack types might increase the bias of detection approaches toward specific situations. To limit the problem of model bias and enable accurate evaluation of detection algorithms, it is therefore crucial to consider datasets that are as free as possible from internal biases while also being sufficiently representative of real-world data. However, dataset bias has been considered mainly for the benchmark datasets used in the field of computer vision [138], [139], whereas there is no analysis in the literature of the bias of benchmark datasets for intrusion detection. Therefore, an evaluation of dataset bias for IDSs may contribute to a fairer assessment of the various algorithms that have been proposed in the field of cybersecurity.

In addition to dataset bias, a few works in the literature address other issues related to public benchmark datasets, such as repeated data, missing values, incorrect labeling [137], or an optimistic number of false alarms due to considering specific situations in a nonrealistic way [140].

B. Novel Features

As the number of methodologies that are able to achieve high accuracy on known datasets increases, attack patterns tend to evolve to better cheat the existing IDSs. This evolution, which can arise in nonstationary environments, is known as concept shift and occurs as the definitions of attacks change over time [141].

For instance, the work presented in [142] shows that some low-frequency DDoS attacks that appear in newer datasets exhibit a higher degree of similarity to normal data traffic than do similar attacks in older datasets. As a consequence, in recent cases, some features are less effective in detecting such attacks than they are in detecting older attack patterns.

Therefore, it remains an open research issue to investigate whether the available features in known benchmark datasets are sufficient to achieve high detection rates even in the presence of changing attack patterns or whether it will be necessary to add new features to maintain a high level of detection accuracy.

V. CONCLUSION

In this review, we have analyzed ML-based approaches to cybersecurity and IDSs, with a specific focus on the most recent

methods based on DL, which represent the current state of the art for intrusion detection in network traffic. Specifically, we have considered methods based on DBNs, AEs, CNNs, LSTM networks, and GANs. In contrast to previous surveys, this review considers studies that use common benchmark datasets to ensure a fair evaluation and comparison of the proposed algorithms.

To provide a reference for how recent cybersecurity methods use benchmark datasets for intrusion detection, in this survey, we have also reviewed the main datasets used for this purpose by highlighting their potential for training effective ML-based algorithms. In particular, we have considered the data collection procedures, the distributions of feature and attack types, and dataset reliability criteria.

By providing a survey of ML and DL approaches, along with descriptions of the benchmark datasets considered when developing recent methods, this review aims to provide a practical road map for researchers in academia and industry working in the field of ML and DL for cybersecurity applications.

REFERENCES

- [1] S. Muggleton, "Alan Turing and the development of artificial intelligence," *AI Commun.*, vol. 27, no. 1, pp. 3–10, 2014.
- [2] "WannaCry ransomware attack," 2017. [Online]. Available: https://en.wikipedia.org/wiki/WannaCry_ransomware_attack
- [3] "Hacked consumers don't forgive companies who lose their data. Bad news for Yahoo," 2016. [Online]. Available: <https://www.thestreet.com/personal-finance/how-long-if-ever-will-it-take-consumers-to-trust-yahoo-after-the-data-breach-13752314>
- [4] McAfee, "McAfee labs threats report," 2019. [Online]. Available: <https://www.mcafee.com/enterprise/en-us/assets/reports/rp-quarterly-threats-aug-2019.pdf>
- [5] R. Bhadoria, "Security architecture for cloud computing," in *Cyber Security and Threats: Concepts, Methodologies, Tools, and Applications*. Hershey, PA, USA: IGI Global, 2018, pp. 729–755.
- [6] M. Swarnkar and R. Bhadoria, "Security aspects in utility computing," in *Emerging Research Surrounding Power Consumption and Performance Issues in Utility Computing*. Hershey, PA, USA: IGI Global, 2016, pp. 262–275.
- [7] S. Dorbala and R. Bhadoria, "Analysis for security attacks in cyber-physical systems," in *Cyber-Physical Systems: A Computational Perspective*. London, U.K.: Chapman & Hall, 2015, pp. 395–414.
- [8] S. K. Khaitan and J. D. McCalley, "Design techniques and applications of cyberphysical systems: A survey," *IEEE Syst. J.*, vol. 9, no. 2, pp. 350–365, Jun. 2015.
- [9] R. Sandhu and P. Samarati, "Authentication, access control and intrusion detection," in *CRC Handbook of Computer Science and Engineering*. Boca Raton, FL, USA: CRC Press, 1997, pp. 1929–1948.
- [10] S. Han, M. Xie, H. Chen, and Y. Ling, "Intrusion detection in cyber-physical systems: Techniques and challenges," *IEEE Syst. J.*, vol. 8, no. 4, pp. 1052–1062, Dec. 2014.
- [11] T. T. T. Nguyen and G. Armitage, "A survey of techniques for internet traffic classification using machine learning," *IEEE Commun. Surv. Tut.*, vol. 10, no. 4, pp. 56–76, Oct.–Dec. 2008.
- [12] R. Donida Labati, A. Genovese, V. Piuri, F. Scotti, and S. Vishwakarma, "Computational intelligence in cloud computing," in *Recent Advances in Intelligent Engineering: Volume Dedicated to Imre J. Rudas' Seventieth Birthday*. New York, NY, USA: Springer, 2020, pp. 111–127.
- [13] Y. Cai, A. Genovese, V. Piuri, F. Scotti, and M. Siegel, "IoT-based architectures for sensing and local data processing in ambient intelligence: Research and industrial trends," in *Proc. IEEE Int. Instrum. Meas. Technol. Conf.*, 2019, pp. 1–6.
- [14] Z. M. Fadlullah *et al.*, "State-of-the-art deep learning: Evolving machine intelligence toward tomorrow's intelligent network traffic control systems," *IEEE Commun. Surv. Tut.*, vol. 19, no. 4, pp. 2432–2455, Oct.–Dec. 2017.
- [15] S. Pouyanfar *et al.*, "A survey on deep learning: Algorithms, techniques, and applications," *ACM Comput. Surv.*, vol. 51, no. 5, 2018, Art. no. 92.

- [16] M. H. Bhuyan, D. K. Bhattacharyya, and J. K. Kalita, "Network anomaly detection: Methods, systems and tools," *IEEE Commun. Surv. Tut.*, vol. 16, no. 1, pp. 303–336, Jan.–Mar. 2014.
- [17] A. L. Buczak and E. Guven, "A survey of data mining and machine learning methods for cyber security intrusion detection," *IEEE Commun. Surv. Tut.*, vol. 18, no. 2, pp. 1153–1176, Apr.–Jun. 2016.
- [18] P. Mishra, V. Varadharajan, U. Tupakula, and E. S. Pilli, "A detailed investigation and analysis of using machine learning techniques for intrusion detection," *IEEE Commun. Surv. Tut.*, vol. 21, no. 1, pp. 686–728, Jan.–Mar. 2019.
- [19] D. Kwon, H. Kim, J. Kim, S. C. Suh, I. Kim, and K. J. Kim, "A survey of deep learning-based network anomaly detection," *Cluster Comput.*, vol. 22, pp. 949–961, 2017.
- [20] E. Hodo, X. Bellekens, A. Hamilton, C. Tachtatzis, and R. Atkinson, "Shallow and deep networks intrusion detection system: A taxonomy and survey," Ithaca, NY, USA, 2017, pp. 1–43.
- [21] Y. Xin *et al.*, "Machine learning and deep learning methods for cybersecurity," *IEEE Access*, vol. 6, pp. 35365–35381, 2018.
- [22] H. Hindy *et al.*, "A taxonomy and survey of intrusion detection system design techniques, network threats and datasets," *CoRR*, vol. abs/1806.03517, pp. 1–35, 2018.
- [23] A. Praseed and P. S. Thilagam, "DDoS attacks at the application layer: Challenges and research perspectives for safeguarding web applications," *IEEE Commun. Surv. Tut.*, vol. 21, no. 1, pp. 661–685, Jan.–Mar. 2019.
- [24] B. B. Zarpelo, R. S. Miani, C. T. Kawakani, and S. C. de Alvarenga, "A survey of intrusion detection in Internet of Things," *J. Netw. Comput. Appl.*, vol. 84, no. C, pp. 25–37, 2017.
- [25] C. Tsai, C. Lai, M. Chiang, and L. T. Yang, "Data mining for Internet of Things: A survey," *IEEE Commun. Surv. Tut.*, vol. 16, no. 1, pp. 77–97, Jan.–Mar. 2014.
- [26] A. Nisioti, A. Mylonas, P. D. Yoo, and V. Katos, "From intrusion detection to attacker attribution: A comprehensive survey of unsupervised methods," *IEEE Commun. Surv. Tut.*, vol. 20, no. 4, pp. 3369–3388, Oct.–Dec. 2018.
- [27] R. Abdulhammed, M. Faezipour, and K. M. Elleithy, "Network intrusion detection using hardware techniques: A review," in *Proc. IEEE Long Island Syst., Appl. Technol. Conf.*, 2016, pp. 1–7.
- [28] J. Kim, P. J. Bentley, U. Aickelin, J. Greensmith, G. Tedesco, and J. Twycross, "Immune system approaches to intrusion detection—A review," *Natural Comput.*, vol. 6, no. 4, pp. 413–466, 2007.
- [29] A. Volkova, M. Niedermeier, R. Basmadjian, and H. de Meer, "Security challenges in control network protocols: A survey," *IEEE Commun. Surv. Tut.*, vol. 21, no. 1, pp. 619–639, Jan.–Mar. 2019.
- [30] O. Savas and J. Deng, *Big Data Analytics in Cybersecurity*. New York, NY, USA: Auerbach, 2017.
- [31] F. Pacheco, E. Exposito, M. Gineste, C. Baudoin, and J. Aguilar, "Towards the deployment of machine learning solutions in network traffic classification: A systematic survey," *IEEE Commun. Surv. Tut.*, vol. 21, no. 2, pp. 1988–2014, Apr.–Jun. 2019.
- [32] "LibPCAP," 2019. [Online]. Available: <https://www.tcpdump.org>
- [33] "WinPCAP," 2018. [Online]. Available: <https://www.winpcap.org>
- [34] "Snort," 2020. [Online]. Available: <https://www.snort.org>
- [35] "Wireshark," 2020. [Online]. Available: <https://www.wireshark.org>
- [36] "TShark," 2020. [Online]. Available: <https://www.wireshark.org/docs/man-pages/tshark.html>
- [37] "TCPDump," 2019. [Online]. Available: <https://www.tcpdump.org>
- [38] "Networkminer," 2019. [Online]. Available: <https://www.netresec.com/?page=NetworkMiner>
- [39] "Rapidminer," 2020. [Online]. Available: <https://rapidminer.com>
- [40] "Scapy," 2019. [Online]. Available: <https://scapy.net>
- [41] "Cisco Netflow," 2012. [Online]. Available: <https://www.cisco.com/c/en/us/products/ios-nx-os-software/ios-netflow/index.html>
- [42] "Nfdump," 2020. [Online]. Available: <https://github.com/paag/nfdump>
- [43] D. Jurafsky and J. H. Martin, *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, 1st ed. Englewood Cliffs, NJ, USA: Prentice-Hall, 2000.
- [44] R. Hofstede *et al.*, "Flow monitoring explained: From packet capture to data analysis with NetFlow and IPFIX," *IEEE Commun. Surv. Tut.*, vol. 16, no. 4, pp. 2037–2064, Oct.–Dec. 2014.
- [45] M. Tavallaei, E. Bagheri, W. Lu, and A. A. Ghorbani, "A detailed analysis of the KDD CUP 99 data set," in *Proc. IEEE Symp. Comput. Intell. Secur. Defense Appl.*, 2009, pp. 1–6.
- [46] X. Jing, Z. Yan, and W. Pedrycz, "Security data collection and data analytics in the internet: A survey," *IEEE Commun. Surv. Tut.*, vol. 21, no. 1, pp. 586–618, Jan.–Mar. 2019.
- [47] J. Fonseca, M. Vieira, and H. Madeira, "Testing and comparing web vulnerability scanning tools for SQL injection and XSS attacks," in *Proc. Pac. Rim Int. Symp. Dependable Comput.*, 2007, pp. 365–372.
- [48] A. Lazarevic, V. Kumar, and J. Srivastava, "Intrusion detection: A survey," *Managing Cyber Threats*, vol. 5, pp. 19–78, 2005.
- [49] "AWID2018 dataset," University of the Aegean, 2018. [Online]. Available: <http://icsdweb.aegean.gr/awid/features.html>
- [50] "Intrusion Detection Evaluation Dataset," Canadian Institute for Cybersecurity, 2017. [Online]. Available: <https://www.unb.ca/cic/datasets/ids-2017.html>
- [51] "KDD Cup 1999," University of California, Irvine, 1999. [Online]. Available: <http://www.kdd.org/kdd-cup/view/kdd-cup-1999>
- [52] "Traffic data from Kyoto University's Honeypots," Kyoto University, 2015. [Online]. Available: http://www.takakura.com/Kyoto_data
- [53] N. Moustafa and J. Slay, "UNSW-NB15: A comprehensive data set for network intrusion detection systems," in *Proc. Mil. Commun. Inf. Syst. Conf.*, 2015, pp. 1–6.
- [54] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "Toward generating a new intrusion detection dataset and intrusion traffic characterization," in *Proc. Int. Conf. Inf. Syst. Secur. Privacy*, 2018, pp. 108–116.
- [55] "1998 DARPA intrusion detection evaluation dataset," Massachusetts Institute of Technology, 1998. [Online]. Available: <https://www.ll.mit.edu/r-d/datasets/1998-darpa-intrusion-detection-evaluation-dataset>
- [56] B. Sangster *et al.*, "Toward instrumenting network warfare competitions to generate labeled datasets," in *Proc. Cyber Secur. Exp. Test*, 2009, p. 9.
- [57] "Intrusion detection evaluation dataset (ISCXIDS2012)," Canadian Institute for Cybersecurity, 2012. [Online]. Available: <https://www.unb.ca/cic/datasets/ids.html>
- [58] A. Shiravi, H. Shiravi, M. Tavallaei, and A. A. Ghorbani, "Toward developing a systematic approach to generate benchmark datasets for intrusion detection," *Comput. Secur.*, vol. 31, no. 3, pp. 357–374, 2012.
- [59] "DoS dataset (CIC DoS dataset 2017)," Canadian Institute for Cybersecurity, 2017. [Online]. Available: <https://www.unb.ca/cic/datasets/dos-dataset.html>
- [60] ALDAPA, "Gure-KDDCUP dataset," 2008. [Online]. Available: <http://www.sc.ehu.es/acwaldap/gureKddcup>
- [61] I. Perona, I. Gurrutxaga, O. Arbelaitz, J. I. Martín, J. Muguerza, and J. M. Pérez, "Service-independent payload analysis to improve intrusion detection in network traffic," in *Proc. 7th Australas. Data Mining Conf.*, 2008, pp. 171–178.
- [62] "Cyber Defense Exercise (CDX)," National Security Agency, 2001. [Online]. Available: <https://apps.nsa.gov/iaarchive/programs/cyber-defense-exercise/index.cfm>
- [63] I. Homoliak, M. Barabas, P. Chmelar, M. Drozd, and P. Hanacek, "ASN: Advanced security network metrics for attack vector description," in *Proc. Secur. Manage.*, 2013, pp. 350–358.
- [64] R. Pang, M. Allman, V. Paxson, and J. Lee, "The devil and packet trace anonymization," *SIGCOMM Comput. Commun. Rev.*, vol. 36, no. 1, pp. 29–38, 2006.
- [65] S. Saad *et al.*, "Detecting P2P botnets through network behavior analysis and machine learning," in *Proc. Int. Conf. Privacy, Secur. Trust*, 2011, pp. 174–180.
- [66] R. Fontugne, P. Borgnat, P. Abry, and K. Fukuda, "MAWILab: Combining diverse anomaly detectors for automated anomaly labeling and performance benchmarking," in *Proc. Conf. Emerg. Netw. Exp. Technol.*, 2010, pp. 1–12.
- [67] S. García, M. Grill, J. Stiborek, and A. Zunino, "An empirical comparison of botnet detection methods," *Comput. Secur.*, vol. 45, pp. 100–123, 2014.
- [68] "UMassTraceRepository," Laboratory for Advanced Software Systems, University of Massachusetts Amherst, 2018. [Online]. Available: <http://traces.cs.umass.edu/index.php/Network/Network>
- [69] A. Sperotto, R. Sadre, F. van Vliet, and A. Pras, "A labeled data set for flow-based intrusion detection," in *IP Operations and Management* (Lecture Notes in Computer Science). New York, NY, USA: Springer, 2009, pp. 39–50.
- [70] "Data collection, curation and sharing," Center for Applied Internet Data Analysis, 2018. [Online]. Available: <https://www.caيدا.org/data/>
- [71] G. Creech and J. Hu, "Generation of a new IDS test dataset: Time to retire the KDD collection," in *Proc. IEEE Wireless Commun. Netw. Conf.*, 2013, pp. 4487–4492.

- [72] R. Salakhutdinov and G. Hinton, "Deep Boltzmann machines," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2009, pp. 448–455.
- [73] I. Goodfellow, Y. Bengio, and A. Courville, "Autoencoders," in *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.
- [74] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [75] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [76] I. J. Goodfellow *et al.*, "Generative adversarial nets," in *Proc. Adv. Neural Inf. Process. Syst.*, 2014, pp. 2672–2680.
- [77] G. E. Hinton, "Training products of experts by minimizing contrastive divergence," *Neural Comput.*, vol. 14, no. 8, pp. 1771–1800, 2002.
- [78] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural Comput.*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [79] M. A. Salama, H. Eid, R. Ramadan, A. Darwish, and A. E. Hassanien, "Hybrid intelligent intrusion detection scheme," *Adv. Intell. Soft Comput.*, vol. 96, pp. 295–302, 2011.
- [80] G. Zhao, C. Zhang, and L. Zheng, "Intrusion detection using deep belief network and probabilistic neural network," in *Proc. IEEE Int. Conf. Comput. Sci. Eng.*, 2017, vol. 1, pp. 639–642.
- [81] N. Gao, L. Gao, Q. Gao, and H. Wang, "An intrusion detection model based on deep belief networks," in *Proc. Int. Conf. Adv. Cloud Big Data*, 2014, pp. 247–252.
- [82] M. Z. Alom, V. Bontupalli, and T. M. Taha, "Intrusion detection using deep belief networks," in *Proc. Nat. Aerosp. Electron. Conf.*, 2015, pp. 339–344.
- [83] K. Alrawashdeh and C. Purdy, "Toward an online anomaly intrusion detection system based on deep learning," in *Proc. IEEE Int. Conf. Mach. Learn. Appl.*, 2016, pp. 195–200.
- [84] E. R. Merino, F. M. Castrillejo, J. D. Pin, and D. B. Prats, "Weighted contrastive divergence," *Neural Netw.*, vol. 114, pp. 147–156, 2019.
- [85] NVIDIA, "CUDA," 2020. [Online]. Available: <https://developer.nvidia.com/cuda-zone>
- [86] B. Abolhasanzadeh, "Nonlinear dimensionality reduction for intrusion detection using auto-encoder bottleneck features," in *Proc. 7th Conf. Inf. Knowl. Technol.*, 2015, pp. 1–5.
- [87] M. Yousefi-Azar, V. Varadharajan, L. Hamey, and U. Tupakula, "Autoencoder-based feature learning for cyber security applications," in *Proc. Int. Joint Conf. Neural Netw.*, 2017, pp. 3854–3861.
- [88] V. L. Cao, M. Nicolau, and J. McDermott, "A hybrid autoencoder and density estimation model for anomaly detection," in *Proc. Int. Conf. Parallel Problem Solving Nature*, 2016, pp. 717–726.
- [89] B. Zong *et al.*, "Deep autoencoding Gaussian mixture model for unsupervised anomaly detection," in *Proc. Int. Conf. Learn. Represent.*, 2018, pp. 1–19.
- [90] A. Y. Javaid, Q. Niyaz, W. Sun, and M. Alam, "A deep learning approach for network intrusion detection system," in *Proc. Int. Conf. Bio-Inspired Inf. Commun. Technol.*, 2016, pp. 21–26.
- [91] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *J. Mach. Learn. Res.*, vol. 11, pp. 3371–3408, 2010.
- [92] Y. Yu, J. Long, and Z. Cai, "Network intrusion detection through stacking dilated convolutional autoencoders," *Secur. Commun. Netw.*, vol. 2017, pp. 1–10, 2017.
- [93] Q. Niyaz, W. Sun, and A. Y. Javaid, "A deep learning based DDoS detection system in software-defined networking (SDN)," *EAI Endorsed Trans. Secur. Saf.*, vol. 4, no. 12, pp. 1–10, 2017.
- [94] N. Shone, T. N. Ngoc, V. D. Phai, and Q. Shi, "A deep learning approach to network intrusion detection," *IEEE Trans. Emerg. Topics Comput. Intell.*, vol. 2, no. 1, pp. 41–50, Feb. 2018.
- [95] F. Farahnakian and J. Heikkonen, "A deep auto-encoder based approach for intrusion detection system," in *Proc. Int. Conf. Adv. Commun. Technol.*, 2018, pp. 178–183.
- [96] L. R. Parker, P. D. Yoo, T. A. Asyhari, L. Chermak, Y. Jhi, and K. Taha, "DEMISE: Interpretable deep extraction and mutual information selection techniques for IoT intrusion detection," in *Proc. Int. Conf. Availability, Rel. Secur.*, 2019, pp. 98:1–98:10.
- [97] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," in *Proc. Int. Conf. Learn. Represent.*, 2014, pp. 1–14.
- [98] Q. P. Nguyen, K. W. Lim, D. M. Divakaran, K. H. Low, and M. C. Chan, "GEE: A gradient-based explainable variational autoencoder for network anomaly detection," in *Proc. IEEE Conf. Commun. Netw. Secur.*, 2019, pp. 91–99.
- [99] G. Maciá-Fernández, J. Camacho, R. Magán-Carrión, P. García-Teodoro, and R. Thérón, "UGR'16: A new dataset for the evaluation of cyclostationarity-based network IDSs," *Comput. Secur.*, vol. 73, pp. 411–424, 2018.
- [100] L. Vu, V. L. Cao, Q. U. Nguyen, D. N. Nguyen, D. T. Hoang, and E. Dutkiewicz, "Learning latent distribution for distinguishing network traffic in intrusion detection system," in *Proc. IEEE Int. Conf. Commun.*, 2019, pp. 1–6.
- [101] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Commun. ACM*, vol. 60, no. 6, pp. 84–90, 2017.
- [102] A. Genovese, V. Piuri, K. N. Plataniotis, and F. Scotti, "PalmNet: Gabor-PCA convolutional networks for touchless palmprint recognition," *IEEE Trans. Inf. Forensics Secur.*, vol. 14, no. 12, pp. 3160–3174, Dec. 2019.
- [103] R. Donida Labati, A. Genovese, E. Muñoz, V. Piuri, and F. Scotti, "A novel pore extraction method for heterogeneous fingerprint images using convolutional neural networks," *Pattern Recognit. Lett.*, vol. 113, no. 1, pp. 58–66, 2018.
- [104] A. Genovese, V. Piuri, F. Scotti, and S. Vishwakarma, "Touchless palmprint and finger texture recognition: A deep learning fusion approach," in *Proc. IEEE Int. Conf. Comput. Intell. Virtual Environ. Meas. Syst. Appl.*, 2019, pp. 1–6.
- [105] R. Donida Labati, A. Genovese, E. Muñoz, V. Piuri, and F. Scotti, "Applications of computational intelligence in industrial and environmental scenarios," in *Learning Systems: From Theory to Practice*, vol. 756. New York, NY, USA: Springer, 2018, pp. 29–46.
- [106] Z. Li, Z. Qin, K. Huang, X. Yang, and S. Ye, "Intrusion detection using convolutional neural networks for representation learning," in *Neural Information Processing*. New York, NY, USA: Springer, 2017, pp. 858–866.
- [107] "One-hot encoding," 2020. [Online]. Available: <https://www.sciencedirect.com/topics/computer-science/one-hot-encoding>
- [108] M. Kalash, M. Rochan, N. Mohammed, N. D. B. Bruce, Y. Wang, and F. Iqbal, "Malware classification with deep convolutional neural networks," in *Proc. IFIP Int. Conf. New Technol., Mobility Secur.*, 2018, pp. 1–5.
- [109] T. Kim, S. C. Suh, H. Kim, J. Kim, and J. Kim, "An encoding technique for CNN-based network anomaly detection," in *Proc. Big Data*, 2018, pp. 2960–2965.
- [110] R. Blanco, P. Malagón, J. J. Cilla, and J. M. Moya, "Multiclass network attack classifier using CNN tuned with genetic algorithms," in *Proc. Int. Symp. Power Timing Model., Optim. Simul.*, 2018, pp. 177–182.
- [111] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. Comput. Vision Pattern Recognit.*, 2016, pp. 770–778.
- [112] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proc. Comput. Vision Pattern Recognit.*, 2015, pp. 1–9.
- [113] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. Int. Conf. Learn. Represent.*, 2015, pp. 1–14.
- [114] K. Wu, Z. Chen, and W. Li, "A novel intrusion detection model for a massive network using convolutional neural networks," *IEEE Access*, vol. 6, pp. 50850–50859, 2018.
- [115] U. Çekmez, Z. Erdem, A. G. Yavuz, O. K. Sahingoz, and A. Buldu, "Network anomaly detection with deep learning," in *Proc. 26th Signal Process. Commun. Appl. Conf.*, 2018, pp. 1–4.
- [116] M. Ito and H. Iyatomi, "Web application firewall using character-level convolutional neural network," in *Proc. Int. Colloq. Signal Process. Appl.*, 2018, pp. 103–106.
- [117] S. Z. Lin, Y. Shi, and Z. Xue, "Character-level intrusion detection based on convolutional neural networks," in *Proc. Int. Joint Conf. Neural Netw.*, 2018, pp. 1–8.
- [118] Y. Xiao, C. Xing, T. Zhang, and Z. Zhao, "An intrusion detection model based on feature reduction and convolutional neural networks," *IEEE Access*, vol. 7, pp. 42210–42219, 2019.
- [119] G. Feng, B. Li, M. Yang, and Z. Yan, "V-CNN: Data visualizing based convolutional neural network," in *Proc. IEEE Int. Conf. Signal Process., Commun. Comput.*, 2018, pp. 1–6.
- [120] R. Ronen, M. Radu, C. Feuerstein, E. Yom-Tov, and M. Ahmadi, "Microsoft malware classification challenge," 2018, pp. 1–7, *arXiv:abs/1802.10135*.
- [121] S.-N. Nguyen, V.-Q. Nguyen, J. Choi, and K. Kim, "Design and implementation of intrusion detection system using convolutional neural network for DoS detection," in *Proc. Int. Conf. Mach. Learn. Soft Comput.*, 2018, pp. 34–38.

- [122] S. Park, M. Kim, and S. Lee, "Anomaly detection for HTTP using convolutional autoencoders," *IEEE Access*, vol. 6, pp. 70884–70901, 2018.
- [123] R. Kruse *et al.*, *Computational Intelligence: A Methodological Introduction*, 2nd ed. New York, NY, USA: Springer, 2016.
- [124] A. Brown, A. Tuor, B. Hutchinson, and N. Nichols, "Recurrent neural network attention mechanisms for interpretable system log anomaly detection," in *Proc. Mach. Learn. Comput. Syst.*, 2018, pp. 1–8.
- [125] G. Kim, H. Yi, J. Lee, Y. Paek, and S. Yoon, "LSTM-based system-call language modeling and robust ensemble method for designing host-based intrusion detection systems," vol. abs/1611.01726, pp. 1–12, 2016.
- [126] F. Jiang *et al.*, "Deep Learning based multi-channel intelligent attack detection for data security," *IEEE Trans. Sustain. Comput.*, pp. 1–11, 2018.
- [127] R. Vinayakumar, K. P. Soman, and P. Poornachandran, "Applying convolutional neural network for network intrusion detection," in *Proc. Int. Conf. Adv. Comput., Commun. Inform.*, 2017, pp. 1222–1228.
- [128] W. Wang *et al.*, "HAST-IDS: Learning hierarchical spatial-temporal features using deep neural networks to improve intrusion detection," *IEEE Access*, vol. 6, pp. 1792–1806, 2018.
- [129] Y. Zhang, X. Chen, L. Jin, X. Wang, and D. Guo, "Network intrusion detection: Based on deep hierarchical network and original flow data," *IEEE Access*, vol. 7, pp. 37004–37016, 2019.
- [130] M. Elbayad, L. Besacier, and J. Verbeek, "Pervasive attention: 2D convolutional neural networks for sequence-to-sequence prediction," in *Proc. 22nd Conf. Computat. Natural Lang. Learn.*, 2018, pp. 1–11.
- [131] A. G. Howard *et al.*, "MobileNets: Efficient convolutional neural networks for mobile vision applications," 2017, pp. 1–9, *arXiv:abs/1704.04861*.
- [132] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," 2018, pp. 1–14, *arXiv:abs/1803.01271*.
- [133] Y. Hong, U. Hwang, J. Yoo, and S. Yoon, "How generative adversarial networks and their variants work: An overview," *ACM Comput. Surv.*, vol. 52, no. 1, pp. 10:1–10:43, 2019.
- [134] A. Genovese, V. Piuri, and F. Scotti, "Towards explainable face aging with generative adversarial networks," in *Proc. IEEE Int. Conf. Image Process.*, 2019, pp. 3806–3810.
- [135] T. Schlegl, P. Seeböck, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs, "Unsupervised anomaly detection with generative adversarial networks to guide marker discovery," in *Information Processing in Medical Imaging*, vol. 10265, Niethammer M. *et al.*, Eds., Cham, Switzerland: Springer, 2017, pp. 146–157.
- [136] H. Zenati, C. S. Foo, B. Lecouat, G. Manek, and V. R. Chandrasekhar, "Efficient GAN-based anomaly detection," 2018, pp. 1–13, *arXiv:abs/1802.06222*.
- [137] A. Gharib, I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "An evaluation framework for intrusion detection dataset," in *Proc. Int. Conf. Inf. Sci. Secur.*, 2016, pp. 1–6.
- [138] A. Torralba and A. A. Efros, "Unbiased look at dataset bias," in *Proc. Comput. Vision Pattern Recognit.*, 2011, pp. 1521–1528.
- [139] T. Tommasi, N. Patricia, B. Caputo, and T. Tuytelaars, "A deeper look at dataset bias," in *Domain Adaptation in Computer Vision Applications*. New York, NY, USA: Springer, 2017, pp. 37–55.
- [140] J. McHugh, "Testing intrusion detection systems: A critique of the 1998 and 1999 DARPA intrusion detection system evaluations as performed by Lincoln Laboratory," *ACM Trans. Inf. Syst. Secur.*, vol. 3, pp. 262–294, 2000.
- [141] J. G. Moreno-Torres, T. Raeder, R. Alaíz-Rodríguez, N. V. Chawla, and F. Herrera, "A unifying view on dataset shift in classification," *Pattern Recognit.*, vol. 45, pp. 521–530, 2012.
- [142] R. F. Fouladi, T. Seifpoor, and E. Anarim, "Frequency characteristics of DoS and DDoS attacks," in *Proc. Signal Process. Commun. Appl. Conf.*, 2013, pp. 1–4.



Dilara Gümüşbaş (Student Member, IEEE) received the M.Sc. degree from the Faculty of Electronics and Computer Science, University of Southampton, Southampton, U.K., in 2014. She is currently working toward the Ph.D. degree with Yıldız Technical University, Istanbul, Turkey.

Her research interests include machine learning, signal, and image processing.



Tulay Yıldırım (Member, IEEE) received the B.Sc. and M.Sc. degrees in electronics and communication engineering from Yıldız Technical University, Istanbul, Turkey, in 1990 and 1992, respectively, and the Ph.D. degree in electrical and electronics engineering from the University of Liverpool, Liverpool, U.K., in 1997.

She is currently a Full Professor with the Department of Electronics and Communications Engineering, Yıldız Technical University. She has authored more than 200 publications in journals, conferences, and books. Her research interests include artificial intelligence, intelligent systems, cyber-physical systems, biomedical instrumentation, biometric person identification and verification systems, electronic circuits and systems, and artificial neural networks.



Angelo Genovese (Member, IEEE) received the Ph.D. degree in computer science from the Università degli Studi di Milano, Milan, Italy, in 2014.

Since 2019, he has been an Assistant Professor in Computer Science with the Università degli Studi di Milano. He has been a Visiting Researcher with the University of Toronto, Toronto, ON, Canada. Original results have been published in more than 40 papers in international journals, proceedings of international conferences, books, and book chapters. His research interests include signal and image processing, three-dimensional reconstruction, artificial intelligence for industrial and environmental monitoring systems, biometric systems, and design methodologies and algorithms for self-adapting systems.

Dr. Genovese is an Associate Editor for the *Journal of Ambient Intelligence and Humanized Computing* (Springer) and *Array* (Elsevier).



Fabio Scotti (Senior Member, IEEE) received the Ph.D. degree in computer engineering from the Politecnico di Milano, Milan, Italy, in 2003.

Since 2015, he has been an Associate Professor in Computer Science with the Università degli Studi di Milano, Milan, Italy. Original results have been published in more than 130 papers in international journals, proceedings of international conferences, books, book chapters, and patents. His research interests include biometric systems, artificial intelligence and machine learning, theory and applications of neural networks, signal and image processing, intelligent measurement systems, environmental and industrial applications, and high-level system design.

Dr. Scotti is an Associate Editor for the *IEEE TRANSACTIONS ON HUMAN-MACHINE SYSTEMS* and *Soft Computing* (Springer). He has been an Associate Editor for the *IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY*.