

# ***D7.4 Publications and updated SotA: State of the Art Review – Encrypted Network Traffic Analysis Solutions***

**Document: WP7-D7.4.1**

**Date: 19 December 2023**

## Contents

Change Log .....	3
<b>Summary.....</b>	<b>4</b>
1.    Introduction.....	5
1.1.    ENTA Solutions .....	6
1.2.    ENTA Solution Development Platform .....	9
2.    SotA of Encrypted Network Application Classification Solutions.....	9
2.1.    Academic Research - Encrypted Network Application Classification .....	10
2.2.    Industry Status.....	13
3.    SotA of IoT Device Security Solutions.....	15
3.1.    IoT Device Discovery.....	15
3.2.    IoT Anomaly and Rogue Device Detection .....	17
3.3.    Industry Status.....	20
4.    Challenges and Trends .....	22
4.1.    Technological Challenges.....	22
4.2.    Technological Trends.....	23
4.3.    Business Trends .....	24
5.    Conclusion .....	26
Appendix: Network Application Activity Identification.....	28
Reference .....	30
Acronym/Glossary .....	36

## Change Log

Version	Submission date	Description of changes	Affected Sections
Initial Draft WP7-D7.4.0	03-23-2023	Initial draft	NA
Updated Version WP7-D7.4.1	12-19-2023	Updated content to reflect advances	All sections

# Summary

This document reviews the State of the Art (SotA) as it pertains to the ENTA (Encrypted Network Traffic Analysis) project. The SotA review focuses on research and advanced technology pertaining to the two use cases supported in the project. The first use case focuses on encrypted network application classification while the second use case focuses on IoT device security, with two sub-cases: (i) IoT device discovery and (ii) IoT anomalies and rogue devices detection.

For each of the use cases, we present associated academic research and industry status. From the review, we note the following:

- The need to deal with encrypted traffic
- The inadequacy of existing network solutions that leverage traditional DPI (Deep Packet Inspection) techniques to classify encrypted network traffic and the resultant need for a solution
- The maturation of academic research on encrypted network analytics
- The extensive reliance on ML/DL based solutions which utilize information about device/traffic characteristics while preserving data privacy
- The need for an encrypted network traffic analytic solution development platform to expedite exploration, evaluation, and deployment of innovative solutions in this domain
- The importance of key technological challenges and industrial trends that will influence future solutions

## 1. Introduction

This document describes the state of the art in the area of encrypted network traffic analysis with focus on solutions which utilize Machine Learning (ML) or Deep Learning (DL) techniques.

Today, more than 80% of Internet traffic is encrypted, with seemingly unabated growth. The introduction of TLS 1.3 with ESNI, QUIC, and the increased use of VPN have significantly reduced visibility into network traffic, rendering traditional techniques for traffic analysis to be ineffective. This lack of visibility into encrypted traffic adversely impacts legitimate uses including for cyber security, security/network operations and Law enforcement. Examples include:

- **SOCs** – Security Operation Centres (SOCs) are unable to detect malware and data exfiltration using encrypted channels, and attack surfaces resulting from rogue IoT device activity
- **IT departments** – IT departments are unable to enforce security and management policy. In addition, provision of quality of service (QoS) is a challenge.
- **LEAs** – Law enforcement agencies (LEA) find it difficult to perform forensics and track illicit data handling activities. Significant cost is incurred as a result.
- **Router/Switch vendors** – Network routers and switches are unable to differentiate traffic types or applications to support Class of Service (COS).
- **Firewall vendors** – next generation firewalls need to monitor and filter traffic based on the detected application or traffic class
- **Military Operations** – Network situational awareness suffers due to limited application awareness

The aforementioned issues motivate the ENTA project to address two specific use cases described in Section 1.1 within the scope of the project timeline. The project is working to develop techniques and solutions to analyse encrypted traffic while preserving user privacy. The two use cases can lead to products/solutions with big market opportunities as described below:

- **Use case-1** – Network Application Classification – The techniques developed will result in complementary technologies which address gaps faced by Deep Packet Inspection (DPI). The ML/DL based solution will co-exist with current DPI technologies. Market reports indicate that the DPI market will grow to \$26B by 2028 from \$5.4B in 2021 with a 25% CAGR.
- **Use case-2** – IoT Device Security – The number of IoT devices has been forecast to grow to 125B by 2030 from 46B in 2022. The ability to discover IoT devices and in particular, detect rogue IoT devices is of key importance as these devices represent a very broad attack surface for malicious threats. We note that the market for IoT device discovery and threat detection for encrypted traffic is a nascent one.

Section 1.1 presents an overview of the two use cases which are the focus of the ENTA project along with a summary of research work in related areas. Section 1.2 describes briefly the software platform infrastructure that expedites the development of ML/DL based encrypted network traffic solutions. Once the ENTA platform is developed, it will be suitable for skilled researchers and practitioners in

this field to utilize the ENTA platform as the basis of research into AI-based network analytics for encrypted traffic.

In the rest of document, Section 2 describes the state of the art (SotA) developments in the area of network application detection while Section 3 describes the SotA for IoT device discovery solutions with focus on discovering rogue IoT devices. Each of the above sections concludes with a discussion of their respective industrial trends. Before the conclusion in Section 5, this document provides a brief discussion of the technological challenges and trends that may affect the SotA of the surveyed solutions in Sections 4.1 and 4.2, respectively. Business trends are described in Section 4.3.

## 1.1. ENTA Solutions

The motivations for the ENTA project define the ENTA objective which is described and scoped here. While the ENTA platform is architected to support a wide range of use-cases which analyse encrypted network traffic, only two specific use cases are developed, implemented and showcased during the duration of the ITEA project – the number of use cases supported are constrained by the ENTA project duration and budget. The ENTA project will address and support research and development for the following two specific use-cases.

- Encrypted Network Application Classification (Use Case-1)
  - The objective of this use case is to classify encrypted traffic applications (e.g., Netflix, Spotify, WhatsApp, etc) and identify their categories (e.g., video stream, audio stream, chat, etc). Although identifying the activities occurring in encrypted traffic applications. (e.g., sending/receiving money/image, watching video, listening to audio, etc) is a natural progression of this use case, it is not part of the current ENTA project. Nevertheless, a survey of this aspect is described in the appendix, “Appendix: Network Application Activity Identification”, for the sake of completeness.
- IoT Device Security (Use Case-2)

This use case consists of the following two subcases:

- **IoT device discovery:** The objective here is to discover all IoT devices connected to a given network, including identifying IoT characteristics and connected non-IoT devices.
- **IoT anomaly and rogue devices detection:** The objective here is to detect and identify misbehaving or unauthorized IoT devices connected to a given network – more specifically, those not operating in normally expected manner especially from the security and operational standpoints such as engaging in attacks and malfunctioning.

The following considerations are important to note before embarking on a more detailed SotA analysis in subsequent sections:

- The scope of applications and IoT devices is broad and numbers in the thousands. Detailed exploration of solutions covering the breadth of available applications/devices is beyond the scope of this SotA review. This document aims not for exhaustiveness but to present key highlights and trends. Thus, existing survey papers will only be briefly described here. We note below a few survey papers of primary relevance.
- **UseCase 1 - Encrypted Application Classification:**
  - Papadogiannaki et al survey encrypted network analysis applications, techniques and countermeasures [Papadogiannaki-2021]. They first described solutions typical for analysis of network traffic before the introduction of encryption and surveyed emerging machine-learning based analytics for classifying protocol and application and identifying application usage of encrypted traffic. Beyond classification, the authors described encrypted network traffic analysis solutions for network security, considering intrusion and malware detection. Finally, they considered how privacy could still be infringed and the existing countermeasures for encrypted traffic analysis.
  - Rezaei and Liu present an overview of commonly available deep learning methods and their application for encrypted traffic classification. The overview of the classification problems starts with classification objective: protocols, application, traffic-types, websites, user actions, operating systems, browsers and others. With respect to data collection, they emphasize reliable labelling, available features, representative dataset, proper dataset pre-processing. For feature categories, they mention time series, header, payload data, and statistical features. For deep learning technique, they list multi-layer perceptron, convolutional neural networks, recurrent neural networks, auto-encoders, and generative adversarial networks. For DL model selection, they caution available input features is a major factor. Finally, they conclude with open problems and opportunities such as the rise of stronger encryption protocols, multi-label classification, middle flow classification, zero-day applications, transfer learning and domain adaptation, and multi-task learning.
- **UseCase 2 - IoT Device Security:**
  - Shen et al survey studies on machine learning-powered encrypted network traffic analysis [Shen-2023]. They first describe the workflow of encrypted traffic analysis with machine learning tools, including traffic collection, traffic representation, encrypted traffic analysis method, and performance evaluation. Then they review existing studies in four areas: network asset identification such as IoT discovery, network characterization including protocol recognition, privacy leakage detection, and anomaly detection. Three of their areas of focus are relevant to ENTA SotA.
  - Chatterjee and Ahmed survey IoT anomaly detection methods and their applications [Chatterjee-2022]. Anomalies are considered to be rare events or observations that

represent a departure from expected behaviour. They consider the following four general aspects of the anomaly detection problem:

- methods such as geometrical, statistical or machine learning approaches
- application intent detection such as “on-going normal activities”, “disruptive activities”, and “outliers”
- anomaly type such as for a single incidence, in a given context, or in a trend
- latency – whether the detection occurs on-line or off-line.

The survey covers a number of IoT-specific topics:

- IoT sensor applications
- Smart-city IoT applications
- Surveillance and video IoT applications
- Network traffic analysis of IoT traffic
- Security applications for IoT devices and infrastructure
- Security applications for IoT data transport

- Jmila et al survey solutions for smart home IoT device classification using machine learning-based network traffic analysis. They analyse the approaches to assess their potential and limitations. They also describe a generic workflow for IoT device classification [Jmila-2022].
- Liu et al survey machine learning (ML) enabling techniques for the detection and identification of IoT devices [Liu-2022]. The ML techniques include learning algorithms, feature engineering for network traffic traces and wireless signals, incremental learning, and abnormality detection. They survey the following aspects of the IoT-device identification and detection problem, without specifically focusing on the issues caused by encrypted network traffic:
  - Device type identification & Device-specific feature identification
  - Device identification based on unsupervised methods
  - Device identification based on DL-based methods such as incremental learning, abnormality detection, hyper parameter, and architecture search
  - Abnormal device detection solutions
- This ENTA SotA focuses on the key aspects of the problems being solved such as classification, discovery, and identification issues of IoT devices or network applications using ML/DL based analysis of encrypted network traffic generated by the device or network applications, including associated issues of performance accuracy and speed
- This ENTA SotA considers solutions in both the academic research and industrial market domains.
- While Sections 2 and 3 present specific SotA pertaining to each of the two use-cases of interest, Section 4 attempts to present generalized trends and challenges which pertain to the general area of ENTA SotA as a whole.

## 1.2. ENTA Solution Development Platform

As the ENTA project is developing a software platform to achieve its objectives, we now present a brief overview of software frameworks which can be utilized for the ENTA project. Specifically, we consider here the need for a tailored software frameworks that can support speedy development and deployment of ENTA solutions.

Specialized software frameworks and platforms help expedite solution development, experimentation, and optimization. Moreover, developing ML/DL-based encrypted traffic analysis solutions as part of the ENTA project, will be benefit greatly from a development platform that is customized and reusable for that class of solutions. Below, we outline the characteristics of such a platform.

Having a platform that is tailored to support network traffic analytics solutions will allow development of high-quality network security solutions with associated life cycle management of the development process. Such a platform leverages open-source components to maximum reuse and minimize reinventing components which already exist.

- There are a few Open-Source Platforms which can be considered as the basis of a specialized platform for encrypted network traffic analytics:
  - HopsWorks ([hopsworks.ai](https://hopsworks.ai)) – Initially developed at KTH University, Sweden, is now managed by Logical Clocks.
  - Prefect ([prefect.io](https://prefect.io)) and MLFlows ([mlflow.org](https://mlflow.org)) – Supported by Databricks. In particular, MLFlows has momentum and a large user community is supporting it. The advancement of MLFlow generic platform technology is worth monitoring.
- There are large number of commercially available generic AI platforms (e.g., Amazon Sagemaker, H2O from Oxdata, Databricks etc.). In general, the advancement of this area is also worth monitoring as well.

We note here that [Chatterjee-2022] surveys briefly a number of frameworks for robust anomaly detection [Zhao-2019, Kayan-2021, Tsogbaatar-2021, An-2020], and for privacy preservation and security [Liu-2021, Cauteruccio-2021, Qureshi-2021].

## 2. SotA of Encrypted Network Application Classification Solutions

As described in Section 1, network traffic classification is a key enabling technology which is required by network planners, network operators, security analysts and LEA (Law Enforcement Agencies) among others. However, with the advent of encryption and its unabated use, existing classification solutions such as deep-packet inspection (DPI), are no longer adequate. In this section, we consider the SotA in the area of encrypted application classification (Section 2.1). We also survey some of the leading-edge advanced technology work as it pertains to industrial developments in this area (Section 2.2).

## 2.1. Academic Research - Encrypted Network Application Classification

In this section, we present the results of a survey into leading academic research in the area of *encrypted* network traffic classification using ML-based approaches in Section 2.1.1 and DL-approaches in Section 2.1.2.

### 2.1.1. Academic Research – ML-based Approaches

ML-based traffic classification techniques have been the focus of increasing studies by researchers in recent years. Among the earliest research work reported in this domain was the study by Moore et al in 2005, which leveraged only header information [Moore-2005]. A subsequent study of ML-based encrypted traffic classification was reported in 2011 [Alshammari-2011] in which the authors evaluated the utility of 3 different ML algorithms to detect encrypted Skype and SSH traffic using binary classification models. In [Alshammari-2015], a ML-based approach is proposed to distinguish encrypted VoIP from other traffic. In [Alan-2016], the authors evaluated three supervised machine methods to identify popular Android apps. In [Khatouni-2021], the authors proposed two ML-based frameworks to classify encrypted traffic service types. Some other ML-based solutions, which are not described in more details later, are [Al-Obaidy-2019, DraperGil-2016, Hajjar-2015, Khatouni-2019, Muehlstein-2011], Taylor-2016, Taylor-2018, Wang-2015, and Zhang-2011].

Below, we present further details regarding some of the related ML-based research referenced in the earlier section.

- **[Moore-2005]** - Moore et al carried out early research investigations into the area of application classification [Moore-2005]. Using a Naïve Bayes estimator, they categorize network traffic by application and are able to obtain 65% accuracy on per-flow classification using header-derived discriminators. With additional refinement, they are able to achieve 95% accuracy which they claim to be better than the traditional classification techniques which yielded 50-70% accuracy. The refinements consist of using kernel density estimation theory and Fast Correlation-Based Filter (FCBF), a method of feature selection and redundancy reduction. The authors' study was carried out on a dataset with 10 traffic categories (e.g., P2P) with each consisting of one or more related applications (e.g., KaZaA, BitTorrent, GnuTella). The study was carried out offline with no research into viability of real-time deployment for the solution. A key benefit of their solution was that it utilized only header information for classification and thus, could classify traffic regardless of whether or not the traffic was encrypted.
- **[Alshammari-2011]** - Continuing the trend of not considering packet content so as to be able to classify encrypted network traffic, Alshammari et al presented a machine learning approach using simple Packet Header feature sets and statistical flow feature sets without using the IP addresses, source/destination ports and payload information to detect encrypted applications in network traffic [Alshammari-2011]. Not only are they able to identify encrypted traffic classes with high accuracy without inspecting payload, IP addresses and port numbers, they are also able to identify which services run in encrypted SSH tunnels. The traffic categories considered are SSH and Skype, discriminating between non-SSH and

SSH and between non-Skype and Skype. Their studies conducted experiments using the AdaBoost, SBB-GP (Symbiotic Bid-based Genetic Programming), and C4.5 algorithms and leveraged either packet header attributes or flow attributes<sup>1</sup>. They found GP to perform better than the two other learning algorithm-based methods. The GP based classifier achieves a range of DR (Discovery Rate) values from 89% to 98% and a range of FDR (False Discovery Rate) values from 0.2% to 0.8%. Moreover, their solutions are able to identify correctly services running over SSH such as interactive login sessions (SHELL), tunnelling (both local and remote), SCP (secure copy), SFTP (secure file transfer) and X11 activities with low false positive rate. Their results also show that the classification-based system trained on data from one network can be employed to run on a different network without new training. For example, for SSH tunnel identification, a subset of University traces is used for training, and the rest of the University traces, public traces (AMP and MAWI) and week 1 and 3 of DARA99 traces are used for testing. Note that the description of these datasets can be found in [Alshammary-2011]. Although differentiating between SSH and non-SSH has been studied, differentiating the content multiplexed inside SSH or VPN tunnel is an open research problem.

- **[Alshammary-2015]** - Alshammary and Nur Zincir-Heywood investigate the robustness<sup>2</sup> of the model's signatures generated by ML-based approaches – C5.0, GP, and AdaBoost – for distinguishing encrypted VoIP (Skype) from other traffic [Alshammary-2015]. The C5.0-based classification approach was found to perform the best on the data sets (University Traces, Italy Traces, NIMS2 Traces, NIMS3 Traces, IPv6 Traces) used. The C5.0-based classifier achieved a 99.6 % DR (Detection Rate) with 0.7 %FPR (False Positive Rate) when trained on one network but tested on another in detecting Skype traffic – a form of solution generalization. With respect to evasion attacks, the signatures generated by the C5.0-based classifier from a statistical feature set and a well-chosen training data set were shown to be robust and not easily evaded. The C5.0-based classifier achieved ≈91 % DR and ≈5 % FPR on the Original Skype flows and ≈85 % DR and ≈5 % FPR on the Altered-Skype flows.
- **[Alan-2016]** - Alan and Kaur evaluated three existing supervised machine learning methods. The first method uses similarity measure (SM) based on the Jaccard's coefficient on traffic bursts which are groups of contiguous incoming or outgoing packets within a TCP connection and are rounded to the nearest 32 bytes. The second method uses Gaussian Naïve Bayes (GNB) classifier on packet sizes of traffic sample with negative values indicating incoming packet sizes. The third method uses Multinomial Naïve Bayes (MNB) classifier using packet sizes as in the second method where term frequency – inverse document frequency transformation and normalization are applied to feature vectors. The proposed solutions assume the launch time traffic characteristics are available, specifically, the packet sizes of apps during their launch time. They first capture network trace of 86,109 app launches by repeatedly running 1,595 applications on 4 distinct Android devices. Using existing supervised learning methods on the packet sizes of the first 64 packets from the same device, they were able to identify popular Android apps with 88% accuracy. When the testing

<sup>1</sup> The flow attributes consist of protocol name, flow duration, and for both traffic flow direction, packet count, byte count, statistical information (min, max, mean, and std) of inter-arrival time, and statistical information of packet length.

<sup>2</sup> Here robustness means effectiveness with respect to traffic traces obtained from different locations/networks, time periods, and padded/morphed traffic.

is on unseen device (but similar operating system/vendor), the accuracy of identifying the apps dropped to 67%.

- **[Khatouni-2021]** - Khatouni et al evaluated two ML-based frameworks (i.e., one-layer classifier and two-layer classifier), undertook feature engineering for optimal features selection, and explored the generality of the solutions in terms of network conditions for the classification of encrypted traffic service types [Khatouni-2021]. Thirteen ML algorithms (Random Forest, Decision Tree, Complement Naïve Bayes, Multinomial Naïve Bayes, k-Nearest Neighbours, Bernoulli Naïve Bayes, Linear Support Vector Machine, Ridge Regression, Nearest Centroid, Support Vector Machine, and Linear Models with Stochastic Gradient Descent) were evaluated on NIMs2018, NIMS2019, and PRI2019 datasets using varying number of extracted service-based and network-based features. They captured and analyzed a large-scale dataset for 9 different services in multiple encrypted channels. Their proposed one-layer framework with the Random Forest model achieved the highest accuracy in identifying the multiple service types, e.g., VoIP, text messaging, video, and audio services, under analysis - without using IP addresses, Port numbers, application header fields, and DPI. Moreover, extensive evaluations showed high accuracy (Average F1 Score: 0.90) is achieved in using the minimum number of features (28 features) and reducing the overfitting effects of the model employed. Finally, the results showed encouraging performances in terms of an additional 5% training data resulting in a robust (well generalized) and portable one-layer model framework that can still perform accurately when tested on a new network in terms of location, time, and traffic volume.

### 2.1.2. Academic Research – DL-based Approaches

Exploiting the high feature learning/extraction capability with minimal feature engineering effort in Deep-Learning approaches has led to more recent DL-based research results such as the one reported DL in [Akbari-2021] that covers various service classes and application classes. Recognizing the need to obviates the need for large labelled datasets, the authors in [Rezaei-2020 and Towhid-2022] propose their respective DL-based solutions. In [Lotfollahi-2020], the authors proposed a deep learning-based solution that integrates both feature extraction and classification phases into one system. Some other DL-based solutions, which are not described in more details later, are [Aceto-2020, Cui-2019, Hou-2019, and Wang-2018].

Below, we present further details regarding some of the related DL-based research referenced in the earlier section.

- **[Lotfollahi-2020]** - Lotfollahi et al propose a deep learning-based approach which integrates both feature extraction and classification phases into one system. Their proposed scheme, called “Deep Packet,” can handle two types of traffic characterization: categorizing network traffic into major classes (e.g., FTP and P2P) and identifying end-user applications therein (e.g., BitTorrent and Skype). Not only can Deep Packet identify encrypted traffic and it can also distinguish between VPN and non-VPN network traffic. The Deep Packet framework uses two deep neural network structures, namely stacked autoencoder (SAE) and convolution neural network (CNN) in order to classify network traffic. Their best result is achieved when CNN is used in Deep Packet as its classification model which achieves a recall value of 0.98 in application identification task and 0.94 in traffic categorization task.

- **[Rezaei-2020]** - Rezaei and Liu propose a semi-supervised approach (using a 1-D CNN model) that obviates the need for large labelled datasets [Rezaei-2020]. They first pre-train a model on a large unlabelled dataset where the input is the time series features of a few sampled packets. Then the learned weights are transferred to a new model that is re-trained on a small labelled dataset. They show that their semi-supervised approach achieves almost the same accuracy as a fully-supervised method with a large labelled dataset, although they use only 20 samples per class. For a dataset of 5 Google services generated from the more challenging QUIC protocol, their approach yields 98% accuracy. To show its efficacy, they also test their approach on two public datasets [Ariel-Dataset-2016, QUIC-Dataset-2018]. Moreover, they study three different sampling techniques and demonstrate that sampling packets from an arbitrary portion of a flow is sufficient for classification.
- **[Akbari-2021]** - Recognizing the need to consider crucial domain-specific features such as traffic shape and timing of packets, Akbari et al developed a neural network architecture based on stacked Long Short-Term Memory (LSTM) layers and Convolutional Neural networks (CNN) [Akbari-2021]. In the peer reviewed publication, the authors claim that their solutions, tested on a real-world mobile traffic dataset from an ISP, achieve an average accuracy of 95% in service classification exclusively over HTTPS and outperform other methods by nearly 50% with fewer false classifications. The eight service classes utilized in their study were: chat, download, games, mail, search, social, streaming, and web. Their Deep Learning (DL) models are generalized to achieve different classification objectives including service-level and application-level classification as well as handle diverse encrypted web protocols such as HTTP/2 and QUIC. Finally, their approach achieves an overall accuracy of 99% on a public dataset for QUIC-based applications. A total of 19 applications were included in the study, with examples such as chat-Facebook, chat-Snapchat, chat-WhatsApp, web-Amazon, web-AppleLocalization, and web-Microsoft among others.
- **[Towhid-2022]** - As in [Rezaei-2020], Towhid and Shariar propose a self-supervised approach (using ResNet34) that can achieve high accuracy on encrypted network traffic classification with a few labelled data [Towhid-2022]. Their method consists of two stages, a pre-training stage to learn with unlabelled data and a fine-tuning stage to piggyback on the weights obtained from the pre-training stage to learn with limited number of labelled data. Their proposed solution evaluated on three publicly available datasets<sup>3</sup> not only achieves high accuracy on encrypted traffic but also has the ability to apply the acquired knowledge on a different dataset. In their experiments, their method outperforms the state-of-the-art baseline methods by 3% in terms of accuracy even with a much lower volume of labelled data

## 2.2. Industry Status

Traditional DPI tools work well on unencrypted network traffic. These tools are used to provide application visibility needed to enable various network services such as traffic management, policy enforcement, QoE (Quality of Experience) assurance, and network security management. With the rise of encrypted network traffic, existing providers of application visibility solutions are now

<sup>3</sup> QUIC, Orange'20 and ISCX VPN-NonVPN datasets are available from [Rezaei-2020], [Akbari-2021] and [Draper-Gil-2016], respectively.

enhancing or complementing their traditional DPI solutions to handle adequately encrypted network traffic. Next, we describe how three companies are addressing the application visibility problem due to encrypted network traffic.

- **[ENEA]** - [ENEA] has been using DPI as a method of filtering data that locates, identifies, and classifies the most relevant datasets to support accurate analyses. Their DPI solution has been enhanced to focus on measuring data that can produce actionable responses in real-time to help generate the metadata needed to feed ML and AI algorithms [ENEA-1]. More specifically, in [ENEA-2], [ENEA] acknowledges that “Traffic classification has also been critical to enforcing policy, optimizing traffic flows, meeting quality of service targets, and generating revenue through differentiated service offerings.” With the rise of encryption, the need to combine machine learning and advanced analytics in network traffic analysis solutions such as that of DPI becomes more urgent.
- **[Sandvine]** - [Sandvine] is aware of how encryption can impact DPI services in needing to shifting the focus from exposing hard visible data to inferring data with a certain degree of uncertainty [Sandvine-1]. So, the ability to recognize encrypted applications, services and traffic is becoming an urgent necessity. [Sandvine] has been doing research and development on using behavioural correlation to reliably link together flows from different protocols and services for identifying mashup applications. Using supervised machine learning models pretrained in-house and validated for accuracy, [Sandvine] develop proprietary models and techniques that are built upon 150 flow parameters. From [Sandvine-2], we note that “Employing these techniques, Sandvine is able not only to broadly classify encrypted traffic into categories (e.g., Web Browsing, Video Streaming, VoIP, etc.) but also to accurately classify unique applications within categories (e.g., Facebook vs. Instagram, WhatsApp vs. Lime, Netflix vs. YouTube) – even when the traffic is encrypted and ESNI is in use.” More information is described in [Sandvine-3].
- **[Rhode & Schwarz Ipoque]** - R&S®PACE 2 and R&S®vPACE are two deep packet inspection (DPI) solutions developed by Ipoque that deliver real-time application visibility. These solutions feature techniques for traffic classification utilize behavioral, statistical and heuristic analyses as well as machine learning (ML) and deep learning (DL). They also come with encrypted traffic intelligence (ETI) and boast advanced features such as first-packet classification, NAT detection or custom service classification that enables implementing one's own DPI signatures by using easy, pre-defined criteria to smoothly extend network visibility on-the-fly [Rhode & Schwarz Ipoque-1]. [Rhode & Schwarz Ipoque] indicated that ETI leverages over 1000 ML and DL features, including statistical, time series, and packet-level features, and the ability to automatically identify and incorporate new features, a form of incremental learning [Rhode & Schwarz Ipoque-2]. More information related to the [Rhode & Schwarz Ipoque] solutions can be found in [Rhode & Schwarz Ipoque-2, Rhode & Schwarz Ipoque-3].

### 3. SotA of IoT Device Security Solutions

To secure IoT networks, operators need the ability to discover their IoT assets and be aware of anomalous activities transpiring in their network. Moreover, being able to pre-emptively identify rogue IoT devices will ensure secure IoT network operation. In this context, we consider the SotA of IoT security from the discovery (Section 3.1), anomaly detection (Section 3.2), and rogue device identification (Section 3.2) point of view. In addition to studying SotA for the academic domain, we review relevant industrial solutions (Section 3.3).

#### 3.1. IoT Device Discovery

In this section, for a given set of IoT devices operating in a network, we highlight solutions that enable quick and accurate discovery of IoT devices by examining encrypted network traffic. The subsection below describes the state-of-the-art from the academic research domain (Section 3.1.1) while SotA in industrial solutions is presented in Section 3.3.

##### 3.1.1. Academic Research

One of the early research studies to detect and classify IoT devices was reported by [Sivanathan-2019]. Their ML based approach achieves 99% identification accuracy for 28 unique IoT devices, 17 of which use *encryption* (TLS/SSL) for communication, with a small dataset of only 50,378 labelled instances. Meidan et al also reported achieving 96% to 99% IoT device type identification depending on various test setup [Meidan-2018]. Later studies with different datasets seem to achieve lower accuracies (96% for [Pashamokhari-2021] and 91% for [Zahid-2022]). Although not always explicitly stated, all these research studies address IoT discovery in the context of encrypted IoT network communication. Further details for each research study are described next.

- **[Meidan-2018]** - Meidan et al apply Random Forest to identify IoT device types. Their dataset included network traffic data from 17 distinct IoT devices, representing nine types of IoT devices from which only eight types are used for training. The devices not of the type used in training are correctly detected as unknown in 96% of the test cases and those that belong to the types used in the training are correctly detected in 99% of the cases. The authors listed the top 10 features found to be important for their solution, although they did not list all the features used. These top 10 features are mainly TCP packet time-to-live statistical characteristics. This explains why better accuracy is achieved as more consecutive sessions of traffic data are analysed.
- **[Sivanathan-2019]** - Sivanathan et al consider an IoT network supporting 28 different IoT devices consisting of cameras, lights, plugs, motion sensors, appliances, and health-monitors [Sivanathan-2019]. After extracting the underlying IoT network traffic characteristics based on the statistical attributes of activity cycles, port numbers, signalling patterns, and cipher suites, they develop a multi-stage machine learning-based classification algorithm that identifies specific IoT devices with over 99% accuracy based on their network activity. More specifically, the 8 key attributes utilized as the basis of features to build the ML models include flow volume, flow duration, average flow rate, device sleep time, server port numbers, DNS queries, NTP queries and cipher suites. Their solution consists of two stages. The first stage, called Stage-0, consists of multiple Naïve Bayes Multinomial classifiers, each

of which takes input a bag of non-overlapping, nominal, and possibly multi-valued network attributes. Specifically, their work identifies respectively 356, 421, and 54 unique words for domain names, remote port numbers, and cipher suite strings, forming three bag types. Furthermore, they combine all corresponding words for non-IoT devices under a column named “others”. The output of each such classifier is a tentative tuple, [Class for the bag, Confidence for the bag]. The second stage, called Stage-1, is a Random Forest classifier which takes as input all outputs of all the Naïve Bayes Multinomial classifier and other network attributes such as flow volume, flow duration, flow rate, sleep time, DNS interval, and NTP interval. The output of Stage-1 is the device identity with a confidence level. The authors collected a total of 50,378 labelled instances captured from different IoT and non-IoT devices generating traffic from either triggered user interactions or autonomously generated activities. The captured instances are randomly split with 70% utilized for training and 30% utilized for testing. The proposed solution achieved a detection accuracy of 99.88%, with a minimal value of RRSE (Root Relative Squared Error) of 5.06 %.

- **[Ou-2019]** - Ou-2019] utilizes the diversity of client-side TLS negotiation time to detect client IoT devices, differentiating between IoT and non-IoT devices. Their evaluation shows that the HTTPS server deployed with their solution, IoTClientDetector, performing ECDHE RSA TLS negotiation with 4096-bit RSA key length can precisely detect client-side IoT devices with true positive rate of around 95% and false positive rate of only 7.8%. The devices included in their study were: Raspberry Pi 3B, Raspberry Pi Zero W, MacBook Pro, and a home-built PC with Windows 10 operating system. The number of obtained non-IoT samples ranges from 61 to 81 and that of IoT samples ranges from 19 to 21 for various encryption modes (RSA 2048, ECDHE 2048, RSA 3072, ECDHE 3072, RSA 4096, and ECDHE 4096). Their solution uses the mean value and standard deviation of  $R_k$  to design rules for detection, where  $R_k$  is obtained as  $R_k = R - R_b$ , where R is the roundtrip time between ServerHelloDone message and ClientKeyExchange message, and  $R_b$  is the roundtrip time between the SYNACK message and ACK message.
- **[Pan-2021]** - [Pan-2021] proposes a One-Class Time Series Meta-Learner called DeepNetPrint that learns the network behavioural fingerprint of IoT devices to identify the presence of IoT activity based on limited availability of network activity traces. Their system utilized an autoencoder and One Class Time Series Prototypical Network components. The model was trained using network traces with traffic from 12 IoT devices and evaluated with network traces of 23 IoT devices (11 of which were withheld and considered unknown IoT devices). In their study, they found that DeepNetPrint was able to identify all 23 IoT devices with accuracy of 81%. Only the following description is provided on the feature set used: “we only extracted the ‘Information’ field containing the high-level information of the network packets that would represent ‘conversational dialogues’ originating from the IoT devices to train the model and infer the devices’ network behavioral fingerprints identification later.”
- **[Pashamokhtari-2021]** - Pashamokhtari et al leverage the results of the analysis of IPFIX telemetry records, which are flow-based data typically collected from the edge of ISP networks [Pashamokhtari-2021]. They analysed three million records emitted, over a period of three months, from a residential testbed with 26 IoT devices, from which 28 flow-level

features are extracted to characterize the network behaviour of these devices. The class of devices included examples such as smart-camera, speaker, door phone and TV to lightbulb, sensor and vacuum cleaner. The features capture data of basic nature such as packet counts, byte count and inter-arrival times using their statistical characteristics such as average value and standard deviation. Moreover, the flow direction, is another feature which is utilized. The authors employ a multi-class ML model leveraging random forest to identify the IoT device types in home networks based on the extracted features from their post-NAT IPFIX records. They utilized 10-fold cross-validation to develop more robust ML models, which resulted in model accuracy of 96% accuracy across all classes. As an additional endeavour unrelated to device discovery, a trust metric is proposed to understand the temporal behaviour of IoT devices – detecting steady/transient functional decline.

- **[Zahid-2022]** - Zahid et al propose using Hierarchical Deep Neural Network (HDNN) to distinguish IoT devices from non-IoT devices based on observed network traffic, achieving an accuracy of 91% [Zahid-2022]. They further classify IoT devices into one of the following six categories: controllers and hubs, cameras, switches and triggers, healthcare devices, electronics, and router, achieving an accuracy of 91.33%. In their testbed, they have 28 IoT and non-IoT devices communicating with one another and collect the network traffic trace for a period of 6 months for their dataset and consider a total of 936,893 samples from the PCAP files for performance evaluation.

Using recursive feature elimination (RFE), they identify 20 features deemed to be most important and optimal. These features have similar characteristics to those of [Pashamokhtari-2021] but also include MAC address and port numbers, which may be biased. Their architecture consists of two DNN networks: the first one oversees identifying whether the device is an IoT device, and the second oversees identifying the class that the detected device belongs to. In Table 6 of [Zahid-2022], they claim their proposed solution surpasses some models that have achieved good results such as those of Random Forest that achieved 82.34% for IoT vs non-IoT device identification<sup>4</sup>, Decision Tree Classifier that achieved 79.88% for IoT vs non-IoT device identification<sup>5</sup>, and CNN-LSTM that achieved 74.5% for IoT category identification as noted in [Bai-2018].

### 3.2. IoT Anomaly and Rogue Device Detection

This section presents the SotA research in the area of IoT Device rogue and anomaly detection. The solutions allow us to determine correctly and quickly whether one or more IoT devices are acting anomalously and/or whether a rogue IoT device has unsanctioned behaviour on a network - this determination is made by applying AI to learn the traffic characteristics of encrypted network traffic

<sup>4</sup> Zahid et al did not describe how or where they obtained the results.

<sup>5</sup> ditto

for the device. The next subsection describes the state-of-the-art in the academic domain (Section 3.2.1). SotA in Industrial solutions are discussed in Section 3.3.

### 3.2.1. Academic Research

Tsogbaatar et al show that IoT anomalies can be detected in order to carry out intelligent flow management in SDN networks. Their optimal solution handles the dataset imbalance and achieves a 99% detection rate for the tested datasets [Tsogbaatar-2021]. Ullah and Mahmoud use a feed-forward neural network to detect anomalous activity in IoT devices, achieving an average accuracy of 97% on the evaluated datasets [Ullah-2022]. Vishwakarma et al demonstrate their solution in real-time on a testbed achieving a detection accuracy ranging from 70% to 97% depending on the datasets tested. Non-IoT devices can also be detected as rogue as in [Ou-2019] which described a non-ML-based solution. Another aspect of rogue identification is identifying the presence of unknown IoT devices based on their inferred network behaviours as considered in [Pan-2021]. Yet another aspect is identifying rogue IoT devices directly as in [Hamza-2019] using ML-methods. Although not explicitly stated, all research studies address IoT anomaly detection in the context of encrypted IoT network communication.

More details for each of the above referenced research studies are described next.

- **[Hamza-2019]** - Leveraging defined MUD (Manufacturer User Description) behavioural profiles, Hamza et al develop machine learning methods for an SDN-based system to identify rogue IoT devices participating in volumetric attacks such as DoS, reflective TCP/UDP/ICMP flooding, and ARP spoofing. As usual, traffic flows are analysed and 20 features are utilized. The features include examples such as total, mean, and standard-deviation of packet/byte count over sliding windows of 2-3- and 4-minutes. The solution detects anomalous patterns of supposedly MUD-compliant network activity via coarse-grained (device-level which is flow direction agnostic) and fine-grained (flow-level which is aware of traffic as having bi-directional flows) SDN telemetry for each IoT device. More specifically, the anomaly detection utilizes the following steps: 1) feature reduction using Principal Component Analysis (PCA); 2) clustering using X-means; and 3) outlier detection using boundary detection and Markov Chain. The solution achieved an accuracy of 89.7% when attempting to detect attacks across all IoT devices. The two datasets and system used in the study are publicly available upon request. The datasets consist of raw packet traces and derived flow counters for data representing both benign and attack traffic generated by 10 IoT devices and collected over a period of one month.
- **[Tsogbaatar-2021]** - Tsogbaatar et al present a deep ensemble learning model-based framework (DeL-IoT) for IoT anomaly detection in SDN controllers to manage traffic flows in SDN switches and IoT devices [Tsogbaatar-2021]. In addition to packet and flow level traffic instances that pass through SDN switches, system metrics of deployed devices and applications are also examined and considered as part of the anomaly detection. In this solution, the SDN controller has a learning module, a detection module, and a flow management module. Features are represented using an auto-encoder and deep feature representation by a non-linear transformation, which are then fed to the learning model

consisting of a stacked auto-encoder. The detection module is a Probabilistic Neural Network (PNN) represented as a Kernel Discriminant Analysis (KDA), a generalization of Linear Discriminant Analysis (LDA), to find the linear combination of features that separate classes. Their solutions are thus of the form of DAE-EPNN (Deep Autoencoder Ensemble Probabilistic Neural Networks) or SAE-EPNN (Stacked Autoencoder Ensemble Probabilistic Neural Networks). They appear to utilize a limited set of raw flow parameters<sup>6</sup> that include duration, protocol, source IP address, destination IP address, source Port, destination Port, packets, bytes, tos, and idle\_age. Although SDN system metrics are also used for anomaly detection, they are not explicitly listed. The flow management module is for the flow control and management at the SDN switch based on the rules of actions as modified in part due to the detected anomalies. Using LSTM, the authors of [Tsogbaatar-2021] also provides IoT device status forecasting, marking each IoT device as either legitimate or anomalous at some specified time point. They use data obtained from the testbed, with datasets including 1% to 9% attack instances for data imbalance scenarios. They also utilize the N-BaloT dataset used in [Meidan-2018]. DeL-IoT achieves a detection rate of 99.8% and 99.9% for their testbed and benchmark datasets, respectively. Moreover, DeL-IoT handles well 1% imbalanced datasets, i.e., datasets with only 1% of attack instances, for both single and multi-class anomalies with  $F_1$  and  $MCC$  (Matthews Correlation Coefficient) measures of around 2-3% better than a single model, PNN.

- **[Ullah-2022]** - Using a feed-forward neural network (FFNN), Ullah and Mahmoud design and develop a system for detecting anomalous activity in IoT devices with a model that uses three types of features: two basic features (protocol and destination port), 11 flow features, and control flag features. The flow features include flow duration, flow bytes/sec, flow inter-arrival mean/standard deviation as well as sub-flow associated information. The 12 control flag features are extracted from TCP packets, including examples such as forward/backward PSH flags, FIN flags, RST flags, and ECE flags. To discern the effectiveness of the trained model, they utilize several datasets to evaluate the ability to identify intrusions and devices that have been compromised. Using all the three types of features, they obtained an average multiclass classification accuracy of 98.66% for all their evaluated datasets (99.86% for BoT-IoT, 93.52% for IoT network intrusion, 99.92 for MQTT-IoT-IDS2020, 99.99 for MQTTset, 99.37 for IoT-23, and 99.28 for IoT-DS2 datasets). They consider both binary classification (normal or attack category) and multiclass classification (normal and different attack categories). These five datasets, accessible at [Ullah-2021], contain traffic representing a total of 18 different types of attacks. This data set includes data from various botnets, where the proposed model can identify the anomaly in various botnets datasets.
- **[Vishwakarma-2022]** - Vishwakarma et al propose a Deep Neural Network-based intrusion detection system (DIDS) to detect IoT network attacks in real-time [Vishwakarma-2022].

<sup>6</sup> The parameters such as packets, bytes, tos, and idle\_age are listed in Algorithm 1 of their paper but are not defined. Moreover, how these parameters are used to derive features is also not described.

Their system's pre-processing stage removes bias features and standardizes the input data. Its trained deep neural network stage detects malicious packets and its final stage, attack identification, identifies packets as either benign or malicious. NF-UQ-NIDS, one of their datasets with 20 different types of networking attacks, is a combination of their other four datasets: NF-UNSWNB15, NF-BoTIoT, NF-ToNLoT, and NF-CSE CICIDS2018 datasets, all of which are converted to a uniform NetFlow format. The combined dataset contains 11,994,893 entries, 9,595,914 of which are used for training and 2,398,979 for testing. The selected features do not rely on packet payload content. They include source port, destination port, protocol, TCP flags, Layer 7 protocol, in byte count, out byte count, in packet count, out packet count. The authors do not mention the number of IoT devices and their types. They compare DIDS to the study presented by [Sarhan-2020] as they also used the same dataset with ensemble algorithms such as Random Forest, Extra Tree, and AdaBoost. In Binary class classification, DIDS achieves the highest accuracy in NF-CSE CIC IDS2018, which is 99.21%. Moreover, in multiclass classification, DIDS achieves the highest accuracy in all the mentioned five datasets, NF-BoT-IoT (83.82%), NF-ToN-IoT (69.53%), NF-CSE CIC IDS2018 (97.21%), NF-UNSW NB15 (97.48%), and NF-UQ NIDS (93.02%). Although real-time intrusion detection was demonstrated on a testbed, the authors did not provide any quantitative performance metrics.

### 3.3. Industry Status

The ultimate goal of the industrial solutions for IoT device security is to ensure that the IoT devices used are legitimate, functioning normally, and not engaged in any malicious activities. Minimizing the threat vector posed by IoT devices requires visibility into those devices including the ability to discover them, identify their characteristics and evaluate their behaviours. In this section, we identify and describe six companies providing solutions in this market segment. Most of them do not provide easily accessible and detailed description of their respective solutions and techniques used – none at all, if any, on the discussion of their solutions when IoT devices communicate with encrypted protocols.

- **[Fortinet]** - Fortinet is a cyber security company that also provides IoT device security solutions. Their Collector agents described in [Fortinet-1] can periodically probe all its nearby neighbouring devices to continuously perform discovery to identify newly connected non-workstation devices in the system, such as printers, cameras, media devices and so on. Their “Inventory Auto Grouping” option enables user to group discovered devices by device type. For example, cameras, network devices, media devices, printers and so on. Basically, Fortinet acknowledges that with the rise of IoT deployment, advanced network security solutions are needed to help network operators to identify every user and device that connects to the network and grant or limit network access appropriately [Fortinet-2]. Note that no detailed description of the above IoT device security solutions is available.
- **[Palo Alto Networks]** - Palo Alto networks is yet another cybersecurity company that also provides IoT device security solution. Their solution uses machine learning techniques to detect vulnerabilities and assess risk based on network traffic behaviours of IoT devices and

dynamically updated threat feeds [Palo Alto Networks-1]. Vulnerability is considered potential when it applies to a specific device type, model, and version number and one or more devices match the specified device type but their model and/or version number are unknown [Palo Alto Networks-2]. Additional information about their solution is available in [Palo Alto Networks-3].

More specifically, Palo Alto Networks offers Zinbox IoT Guardian which performs device discovery, identification, classification, and grouping. From [Palo Alto Networks-4], it “*is an internet of things (IoT) security offering that automates the orchestration of the IoT lifecycle to provide security, management, and optimization of all assets. Zingbox IoT Guardian uses a unique, IoT personality-based approach to secure and manage IoT devices throughout their entire lifecycles, from discovery through retirement. It allows customers to automate threat detection and response for their IT and IoT infrastructures from a single system.*”

- **[Armis+Check Point]** - Armis and Check Point provide visibility and security for managed and unmanaged IoT devices. Without the use of agents or additional hardware, the Armis platform uses the existing network infrastructure to discover and identify every device in any environment—enterprise, medical, industrial, and more. The platform analyses device behaviour to identify risks and threats and provides continuous device risk assessments. Armis discovers devices on and off the network, continuously analyses endpoint behaviour to identify risks and attacks, and protects critical information and systems by identifying suspicious or malicious devices and quarantining them [Armis+Check Point-1].
- **[Axonius]** - From [IoTforAll], Axonius provides a centralized IoT visibility and cybersecurity platform. Using its “*... networking capabilities help monitor industrial controls, mobile devices, cloud systems, including remote and on-premises endpoints. A single device can be used to discover the security coverage gaps of one million devices and 50,000 users.*”. Note that no detailed description of their solutions is publicly available.
- **[Forescout eyeSight]** - From [IoTforAll], “*Forescout eyeSight can discover, classify, and assess a variety of endpoints, including laptops, mobile devices, virtual computers, storage networks, operational technology (OT) systems, and IoT gadgets. It is a powerful, agentless IoT visibility solution that continuously monitors every IP-connected device on a network. Forescout eyeSight has auto-classification capabilities, for it is the world's largest data lake of crowdsourced device intelligence. This data lake offers support for 600 versions of OS, 10,000 device types, including 5,700 vendors and models.*” More information about eyesight is available in [Forescout eyesight-1, Forescout eyesight-2].
- **[Securolytics]** - From [IoTforAll], “*The Securolytics IoT Security platform is a suite of products that helps you secure your internet-connected devices. The suite comprises IoT security, IoT discovery, and IoT control products. Securolytics automates the device discovery capability and identification without requiring agents on endpoints, helping lower the total cost of ownership across organizations.*” More information can be found in [Securolytics-1, Securolytics-2].

## 4. Challenges and Trends

Although we consider two different aspects of security solutions, one with respect to network applications and the other IoT devices, we found that their ML/DL-based solutions share common challenges and trends, and a few specific ones.

Section 4.1 considers the challenges facing current SotA ENTA solutions and Section 4.2 the future trends that may affect them, while Section 4.3 considers the business trends.

### 4.1. Technological Challenges

In this section common challenges or issues are listed first, followed by those affecting specific use cases.

The list of common challenges or issues affecting all, if not most, of the current solutions for the two use cases to be demonstrated in the ENTA project are as follows:

- **Lack of standardized datasets:** Although some datasets are made publicly available, most of them are kept private due to privacy concerns. This hinders comparative study and evaluation.
- **Lack of a development support infrastructure:** The lack of specialized platforms for the ML/DL-based network analytics causes a lot of wasted and duplicate effort, with diversion of focus from model solution development.
- **Lack of generalizability:** The proposed solutions seldom work well with unknown data or data from different network environments. The models often do not handle changes in communication protocols, specific networks infrastructures, user behaviours, or a combination of thereof.
- **Invariant Features:** There is a lack of strategy for selecting invariant features suitable for the dynamically changing problem space.
- **Changing network traffic characteristics:** The premise of ENTA is that network traffic characteristics can uniquely identify applications and IoT devices. However, what happens when the underlying network traffic characteristics behaviour changes? There is a need to monitor and handle changes in network characteristics due hardware and software updates.
- **Unseen applications:** There is a need to adapt solutions dynamically to dynamically changing data characteristics and unknown/unseen applications or security issues.
- **Lack of very large-scale datasets of the order of 100 or more application and IoT classes:** This hinders the opportunity and ability to develop and experiment solutions of such scale.
- **Need for good real-time solutions:** Inspecting fewer packets per flow, though would be faster, cannot guarantee to deliver expected real-time performance.
- **Safeguarding privacy:** There is a need to handle encrypted data without compromising privacy.

- **Handling Countermeasures:** There is a need to handle countermeasures such as leak protection for privacy purpose which render network traffic less differentiable.
- **Lack of explainability of proposed solutions:** Most solutions work for their chosen datasets and fixed test environment. There are inadequate available explanations of how and why they work.

Challenges which specifically affect application classification are related to the coverage of solutions in the following aspects:

- Number of applications
- Type of applications
- Version of application
- Diversity of network architectures and technologies
- Demographic specific factors
- Application and network characteristics that change over time

For the IoT use case, solutions have to contend with the following challenges:

- Different IoT device types
- Different IoT devices
- New and unknown IoT device types and devices
- Number of IoT devices

## 4.2. Technological Trends

We briefly list and discuss technological trends that may affect future ENTA solutions. Some changes in technology are immediate but most of the technology changes occur over a long period of time. These technology changes can be considered in following areas:

- Enabling more effective solutions for the 2 use-cases: Application detection and IoT rogue device and anomaly detection:
  - Advances based on evolving ML/DL techniques – early solutions were ML based but more research efforts have focused on Deep Learning in recent years. Recently, Graph Neural Networks have also been utilized [Shen-2021, Wu-2022].
  - Since the basis of the approach is based on traffic characteristics, advances in new feature to represent characteristics can impact solutions. For example, in recent years, representing an encrypted network traffic byte stream as an image has contributed to effective detection of applications [Shapira-2021, Pathmaperuma-2022].
  - Advances in model development – incremental training and upgrading of models are being reported by researchers, this requires new model creation/update whenever changes in traffic characteristics affecting the original model performance are detected [John-2020].

- Increase in complexity for training Datasets:
  - The proliferation of applications utilizing different encryption methods – TLS 1.3 with ESNI, QUIC, VPN etc [Papadogiannaki-2021].
  - The proliferation of integrated applications e.g., WhatsApp can be used for text chat, voice chat and video chat and image exchanges.
  - Number and type of available IoT devices are increasing significantly for different domains - home, enterprise, industrial, wearables, medical etc. Each type of device has unique characteristics.
- Making target solutions ineffective with advancement of new technologies
  - The rise of quantum communication and computing [Norbert-2021].
  - The rise of active countermeasures against encrypted traffic analytics e.g., techniques adopted in Darknet [Papadogiannaki-2021].
  - IETF work of preserving privacy technologies, (Recent Trends on Privacy-Preserving Technologies under Standardization at the IETF). Some of the examples are DNS over TLS and DNS over QUIC [Dikshit-2023].

### 4.3. Business Trends

According to [IoTWorldToday], the key threats set to emerge in 2023 include the increased importance of AI for both offensive and defensive security. This year was expected to be a record-breaking year for cyber security breach notifications, not only because of the sophistication of threat actors – but also due to larger changes in the world: global unrest, supply chain instability, and soaring inflation. These factors will impact an organization’s ability to mitigate, remediate, or prevent cyber threats. The report also indicates that Ransomware will “flourish,” with this form of malware representing the most prolific and costly kinds seen in recent years.

Below we summarize our analysis of the business trends reviewed as part of this SotA study:

- Application UseCase - At present, Deep Packet Inspection (DPI) technology remains the dominant approach for application detection. There is currently no commercial solution which successfully carries out encrypted application discovery and classification. Vendors such as R&S have an advanced solution but rely mostly on DPI techniques with reliance on work-arounds to address the gap. They have some initial efforts in applying AI to the problem of encrypted application detection. Sandvine has a partial solution for encrypted traffic analysis that's used in the LEA domain.

- IoT UseCase - The market for encrypted IoT device discovery and rogue device detection remains a nascent one and vendors are beginning to work on solutions. Palo Alto has commercialized an IoT device discovery solution. However, this solution is not intended to discovery encrypted IoT devices as it relies on various information sent in the clear in order to carry out its discovery.
- ENTA Platform - while there exist a number of general-purpose AI platforms there are still no special purpose AI platforms which focus on AI-based encrypted network traffic analytics. It is this gap that the ENTA platform seeks to fill. While this ENTA project will develop and show case 2 encrypted network traffic analysis use cases on the ENTA platform, the platform is designed to support rapid and fast paced development of additional use cases.

## 5. Conclusion

This document conducted a survey, assessment and analysis of the State of the Art (SotA) in a specific area of interest pertaining to the ENTA project. In particular, the document focuses on SotA developments in the area of encrypted network traffic classification and analysis. At the outset, the study conducted a brief examination and analysis of suitable open-source platforms as the basis of a specialized platform for encrypted network traffic analytics. As the scope of the current ENTA project focuses on two specific use-cases, we subsequently evaluated SotA in academia and industry as it relates to these use-cases: (i) encrypted network application classification (ii) IoT device security with two sub-cases: (a) IoT device discovery (b) IoT anomaly and rogue device detection. We concluded the study with an analysis of technological challenges/trends and business trends in related areas.

As a result of the SotA study, we arrive at a number of findings, summarized as follows:

- The challenge of handling encrypted traffic continues to loom large with over 90% of network traffic being encrypted today. As the deployment scope and strength of encrypted solutions continue to increase, the challenge of developing visibility solutions becomes ever more pressing.
- The key technological challenges and industrial trends presented and analysed in this study will influence the development of future encrypted network analytic solutions.
- Existing network analytic solutions that leverage traditional DPI (Deep Packet Inspection) techniques remain inadequate to handle encrypted network traffic. Workarounds that had been developed to handle encrypted traffic are beginning to reach their capability limit and new solutions are required.
- Existing industrial solutions for network traffic analysis can continue to be utilized along with whatever workarounds have been developed. However, they require specific complementary solutions to address encrypted traffic. A number of technology vendors in this industry segment are in the early stages of carrying out research and development towards a solution.
- AI-based solutions (ML/DL) remain one of the most promising avenues to provide requisite visibility into encrypted network traffic for the use-cases of the ENTA project as well as other future use-cases. Such solutions preserve the privacy of network content.
- In recent years, there has been a significant increase in the academic research carried out into using AI-based approaches for encrypted network traffic analytics. This research has proven the early viability of machine-learning (ML) approaches in controlled test environments. More recently, deep learning (DL) solutions have been the subject of research in this domain to develop more robust solutions.
- AI-based research in this domain now needs to focus on maturation of proposed solutions to address a number of specific hurdles which remain before the research can be more broadly adapted for use in industry solutions: (i) development of generalized ML/DL models which can operate correctly in diverse network environments – to date, validation of research

results were primarily undertaken with data collected from the same network as the training data (ii) development of AI models which can operate correctly at high speed for deployment in network device data planes – 100Gbps and higher (iii) development of solutions which can operate correctly across the broad range of encryption protocols including TLS, QUIC and VPN-based protocols among others (iv) ability to generate correct prediction results when presented with unseen types of encrypted traffic – traffic types that were not included in the model training.

- Acceleration in maturity of AI-based solutions in different domains is assisted by the existence of collaborative work in the form of open-source libraries, tools, datasets and models. These tools, datasets and models are tailored towards solving a specific problem. A strong example can be observed in the image processing and object detection domain. The domain of encrypted network traffic analytics would benefit immensely from a common platform which will accelerate development of high-quality solutions. This gap, the ENTA project seeks to fill.

## Appendix: Network Application Activity Identification

In this section, for a given set of network applications operating in a network, we highlight solutions that allows us to infer correctly and quickly the application-level activity from *encrypted* network traffic describing the state-of-the-art solutions in the academic research domain.

Different kinds of applications provide different kinds of services and thus, they support different kinds of activities. Here, we consider a few non-exhaustive examples of identifying specific activities occurring within particular applications. The first example, presented in [Aiolli-2019], describes identifying user fund transfer activities on smart-phone-based Bitcoin wallet apps. The second example, described in [Liu-2019], focuses on identifying human behaviours transmitted in encrypted video. The last example is a solution described in [Pathmaperuma-2022] for identifying a total of 92 transpiring activities in 8 applications such as posting various media types, providing different types of comments, and watching different types of media. For their respective datasets, all the studies indicate that they achieve accuracy over 90%. More details for each research study are described next:

- **[Aiolli-2019]** - Considering applications of the same type (smartphone-based Bitcoin wallet apps) and same functionality (sending, receiving, and trading Bitcoin), Aiolli et al used machine learning techniques such as SVM (Support Vector Machines) and RF (Random Forest) to analyse *encrypted network traffic* for those applications and their associated activities. They used statistical features extracted from a sequence of directional packet sizes represented by a sequence of signed integer numbers whose signs are determined by the packet flow direction. The nine types of cryptocurrency applications investigated are: BTC.com, Bitcoin Wallet (Android and iOS), Coinbase, Mycelium, BitPay, Blockchain, Bread, and Copay. The activities considered for these apps are: Open App, Receive Bitcoin, and Send Bitcoin, although many more actions may be available. Their experimental results achieve nearly 95% accuracy in user activity identification for 9 Bitcoin wallet applications, four of which are Android-based with the rest being iOS-based.
- **[Liu-2019]** - Liu et al investigated six typical machine learning algorithms to identify an individual user's daily living behaviours from live video. Behaviour activity that was detected included examples such as watching TV, reading books, styling hair, opening or closing a door, sweeping the ground, getting dressed, drinking and moving around. The authors determined and validated features required to build the requisite model. The features included statistical properties computed from traffic rate change (TRC), gain coefficients in the frequency domain obtained from DFT (Discrete Fourier Transform) of traffic segment, among others [Liu-2019]. The machine learning algorithms investigated include: Naïve Bayes, Logistic Regression, K-Nearest Neighbour, Decision Tree, Gradient Boosting Decision Tree, and Random Forest. Their experimental results show that the user's behaviour captured in *encrypted video traffic* can be identified with 94% accuracy.

- **[Pathmaperuma-2022]** - Considering unknown applications and fine grained in-application activity detection with minimal data, Pathmaperuma et al proposed a Convolutional Neural Network (CNN) in a framework that uses a time window-based approach to split the activity occurring within an *encrypted* traffic flow *into segments*. Their technique considers packet size and time related information. In their solution, these segments are constituted into matrices that serve as input to the CNN model, enabling it to learn to differentiate previously trained (known) and previously untrained (unknown) in-application activities. These in-application activities are then identified via as many such segments as needed. Their approach is able to filter unknown traffic with an average accuracy of 88% and a classification accuracy of 92% once the unknown traffic has been filtered out. The authors claim that their solution yields good results with as little as 0.2s of data exchange for an application. The eight applications considered are Facebook, Instagram, Gmail, Messenger, Skype, Viber, WhatsApp, and YouTube. A total of 92 activities are identified within the 8 applications with the number of distinct activities per application ranging from 5 in Gmail to 22 in Facebook.

## Reference

[Aceto-2020]	Aceto, Giuseppe, Domenico Ciuonzo, Antonio Montieri, and Antonio Pescapé. "Toward Effective Mobile Encrypted Traffic Classification through Deep Learning." <i>Neurocomputing</i> (Amsterdam) 409 (2020): 306–315.
[Aiolli-2019]	Fabio Aiolli, Mauro Conti, Ankit Gangwal, and Mirko Polato. 2019. "Mind your wallet's privacy: identifying Bitcoin wallet apps and user's actions through network traffic analysis," In Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing. ACM, 1484–1491.
[Akbari-2021]	Akbari, et.al., "A Look Behind the Curtain: Traffic Classification in an Increasingly Encrypted Web," <i>ACM on Measurement and Analysis of Computing Systems</i> , Vol 5, 2021
[Al-Obaidy-2019]	F. Al-Obaidy, S. Momtahen, M. F. Hossain and F. Mohammadi, "Encrypted Traffic Classification Based ML for Identifying Different Social Media Applications," 2019 IEEE Canadian Conference of Electrical and Computer Engineering (CCECE), 2019, pp. 1-5, doi: 10.1109/CCECE.2019.8861934.
[Alan-2016]	Hasan Faik Alan and Jasleen Kaur. 2016. "Can Android applications be identified using only TCP/IP headers of their launch time traffic?" In Proceedings of the 9th ACM Conference on Security & Privacy in Wireless and Mobile Networks. ACM, 61–66.
[Alshammari-2011]	Alshammari & Zincir-Heywood, "Can encrypted traffic be identified without port numbers, IP addresses and payload inspection?", <i>J. Computer Networks</i> , Vol 55, 2011
[Alshammari-2015]	Alshammari & Zincir-Heywood, "How Robust Can a Machine Learning Approach Be for Classifying Encrypted VoIP?". <i>J Netw Syst Manage</i> 23, 830–869 (2015). <a href="https://doi.org/10.1007/s10922-014-9324-6">https://doi.org/10.1007/s10922-014-9324-6</a>
[An-2020]	Y. An, F.R. Yu, J. Li, J. Chen, V.C.M. Leung, "Edge intelligence (EI)-enabled HTTP anomaly detection framework for the internet of things (IoT)", <i>IEEE Internet Things J.</i> 8 (5) (2021) 3554–3566, <a href="http://dx.doi.org/10.1109/JIOT.2020.3024645">http://dx.doi.org/10.1109/JIOT.2020.3024645</a> .
[Ariel-Dataset-2016]	A.U.: (2016), <a href="https://drive.google.com/drive/folders/0Bynah7-gERTIdG5UZ2NhNkJMMIk">https://drive.google.com/drive/folders/0Bynah7-gERTIdG5UZ2NhNkJMMIk</a>
[Armis+Check Point]	<a href="https://www.checkpoint.com/downloads/partners/cp-armis-enterprise-iot-security-solution-brief.pdf">https://www.checkpoint.com/downloads/partners/cp-armis-enterprise-iot-security-solution-brief.pdf</a>
[Armis+Check Point-1]	<a href="https://aws.amazon.com/marketplace/pp/prodview-zhbfuevcnjfw">https://aws.amazon.com/marketplace/pp/prodview-zhbfuevcnjfw</a>
[Axonius]	<a href="https://www.axonius.com/">https://www.axonius.com/</a>
[Bai-2018]	L. Bai, L. Yao, S. S. Kanhere, X. Wang, and Z. Yang, "Automatic device classification from network traffic streams of internet of things," in 2018 IEEE 43rd Conference on Local Computer Networks (LCN), Chicago, IL, USA, 2018.
[Cauteruccio-2021]	F. Cauteruccio, L. Cinelli, E. Corradini, G. Terracina, D. Ursino, L. Virgili, C. Savaglio, A. Liotta, G. Fortino, A framework for anomaly detection and classification in multiple IoT scenarios, <i>Future Gener. Comput. Syst.</i> 114 (2021) 322–335, <a href="http://dx.doi.org/10.1016/j.future.2020.08.010">http://dx.doi.org/10.1016/j.future.2020.08.010</a> , URL: <a href="https://www.sciencedirect.com/science/article/pii/S0167739X19335253">https://www.sciencedirect.com/science/article/pii/S0167739X19335253</a> .
[Chatterjee-]	Ayan Chatterjee, Bestoun S. Ahmed, "IoT anomaly detection methods and applications: A

2022]	survey," Internet of Things, Volume 19, 2022, 100568, ISSN 2542-6605, <a href="https://doi.org/10.1016/j.iot.2022.100568">https://doi.org/10.1016/j.iot.2022.100568</a> .
[Cui-2019]	S. Cui, B. Jiang, Z. Cai, Z. Lu, S. Liu, and J. Liu, "A Session-Packets-Based Encrypted Traffic Classification Using Capsule Neural Networks," 2019 IEEE 21st International Conference on High-Performance Computing and Communications; IEEE 17th International Conference on Smart City; IEEE 5th International Conference on Data Science and Systems (HPCC/SmartCity/DSS), Zhangjiajie, China, 2019, pp. 429-436.
[Deforce-2023]	B. Deforce, B. Baesens, J. Diels, and E. Serral Asensio, "Self-Supervised Anomaly Detection of Rogue Soil Moisture Sensors" 2023, <a href="https://doi.org/10.48550/arXiv.2305.05495">https://doi.org/10.48550/arXiv.2305.05495</a>
[Dikshit-2023]	Pratyush Dikshit, Jayasree Sengupta, Vaibhav Bajpai, "Recent Trends on Privacy-Preserving Technologies under Standardization at the IETF," <a href="https://arXiv:2301.01124v1">arXiv:2301.01124v1</a> .
[Draper-Gil-2016]	G. Draper-Gil, A. H. Lashkari, M. S. I. Mamun, and A. A. Ghorbani, "Characterization of Encrypted and VPN Traffic using Time-related Features," in Proceedings of the 2nd International Conference on Information Systems Security and Privacy, Rome, Italy, 2016, pp. 407-414.
[ENEA]	<a href="https://www.enea.com/">https://www.enea.com/</a>
[ENEA-1]	<a href="https://www.enea.com/insights/is-network-dpi-really-a-dead-piece-of-investment/">https://www.enea.com/insights/is-network-dpi-really-a-dead-piece-of-investment/</a>
[ENEA-2]	<a href="https://www.enea.com/insights/yes-encrypted-traffic-can-and-must-be-classified/">https://www.enea.com/insights/yes-encrypted-traffic-can-and-must-be-classified/</a>
[Forescout eyeSight]	<a href="https://www.forescout.com/products/eyesight/">https://www.forescout.com/products/eyesight/</a>
[Forescout eyesight-1]	<a href="https://www.exclusive-networks.com/uk/wp-content/uploads/sites/28/2020/12/UK-VR-Forescout-Data-Sheet-Forescout-eyeSight-1.pdf">https://www.exclusive-networks.com/uk/wp-content/uploads/sites/28/2020/12/UK-VR-Forescout-Data-Sheet-Forescout-eyeSight-1.pdf</a>
[Forescout eyesight-2]	<a href="https://www.forescout.com/resources/forescout-eyesight-datasheet/">https://www.forescout.com/resources/forescout-eyesight-datasheet/</a>
[Fortinet]	<a href="https://www.fortinet.com/">https://www.fortinet.com/</a>
[Fortinet-1]	<a href="https://docs.fortinet.com/document/fortiedr/6.0.0/administration-guide/559754/iot-device-discovery#:~:text=IoT%20device%20discovery%20enables%20you,all%20its%20nearby%20neighboring%20devices">https://docs.fortinet.com/document/fortiedr/6.0.0/administration-guide/559754/iot-device-discovery#:~:text=IoT%20device%20discovery%20enables%20you,all%20its%20nearby%20neighboring%20devices</a>
[Fortinet-2]	<a href="https://www.fortinet.com/demo-center/nac-demo">https://www.fortinet.com/demo-center/nac-demo</a>
[Hajjar-2015]	Amjad Hajjar, Jawad Khalife, Jesús Díaz-Verdejo, "Network traffic application identification based on message size analysis," Journal of Network and Computer Applications, Volume 58, 2015, Pages 130-143, ISSN 1084-8045, <a href="https://doi.org/10.1016/j.jnca.2015.10.003">https://doi.org/10.1016/j.jnca.2015.10.003</a> . ( <a href="https://www.sciencedirect.com/science/article/pii/S1084804515002167">https://www.sciencedirect.com/science/article/pii/S1084804515002167</a> )
[Hamza-2019]	A. Hamza, H. H. Gharakheili, T. Benson, V. Sivaraman, "Detecting Volumetric Attacks on IoT Devices via SDN-Based Monitoring of MUD Activity", ACM SOSR, USA, Apr 2019.
[Hou-2019]	T. Hou, T. Wang, Z. Lu and Y. Liu, "Smart Spying via Deep Learning: Inferring Your Activities from Encrypted Wireless Traffic," 2019 IEEE Global Conference on Signal and Information Processing (GlobalSIP), 2019, pp. 1-5, doi: 10.1109/GlobalSIP45357.2019.8969428.

[IoTforAll]	<a href="https://www.iotforall.com/the-need-for-iot-device-discovery-and-security-in-2022">https://www.iotforall.com/the-need-for-iot-device-discovery-and-security-in-2022</a>
[IoTWorldToday]	<a href="https://www.iotworldtoday.com/security/cybersecurity-threat-predictions-for-the-new-year">https://www.iotworldtoday.com/security/cybersecurity-threat-predictions-for-the-new-year</a>
[Jmila-2022]	H. Jmila, G. Blanc, M. R. Shahid and M. Lazrag, "A Survey of Smart Home IoT Device Classification Using Machine Learning-Based Network Traffic Analysis," in IEEE Access, vol. 10, pp. 97117-97141, 2022, doi: 10.1109/ACCESS.2022.3205023.
[John-2020]	Meenu Mary John, Helena Holmström Olsson, and Jan Bosch. 2020. Developing ML/DL Models: A Design Framework. In Proceedings of the International Conference on Software and System Processes (ICSSP '20). Association for Computing Machinery, New York, NY, USA, 1–10. <a href="https://doi.org/10.1145/3379177.3388892">https://doi.org/10.1145/3379177.3388892</a>
[Khatouni-2019]	<a href="#">A. S. Khatouni and N. Zincir-Heywood, "Integrating Machine Learning with Off-the-Shelf Traffic Flow Features for HTTP/HTTPS Traffic Classification," 2019 IEEE Symposium on Computers and Communications (ISCC), 2019, pp. 1-7, doi: 10.1109/ISCC47284.2019.8969578.</a>
[Khatouni-2021]	<a href="#">Khatouni, Seddigh, Nandy, Zincir-Heywood, "Machine Learning Based Classification Accuracy of Encrypted Service Channels: Analysis of Various Factors". J Netw Syst Manage 29, 8 (2021). https://doi.org/10.1007/s10922-020-09566-5</a>
[Liu-2019]	X. Liu, J. Wang, Y. Yang, Z. Cao, G. Xiong and W. Xia, "Inferring Behaviors via Encrypted Video Surveillance Traffic by Machine Learning," 2019 IEEE 21st International Conference on High Performance Computing and Communications; pp. 273-280.
[Liu-2022]	Y. Liu, J. Wang, J. Li, S. Niu and H. Song, "Machine Learning for the Detection and Identification of Internet of Things Devices: A Survey," in IEEE Internet of Things Journal, vol. 9, no. 1, pp. 298-320, 1 Jan.1, 2022, doi: 10.1109/JIOT.2021.3099028.
[Lotfollahi-2020]	Lotfollahi, M., Jafari Siavoshani, M., Shirali Hossein Zade, R. et al. Deep packet: a novel approach for encrypted traffic classification using deep learning. Soft Comput 24, 1999–2012 (2020). <a href="https://doi.org/10.1007/s00500-019-04030-2">https://doi.org/10.1007/s00500-019-04030-2</a>
[Meidan-2018]	Yair Meidan, Michael Bohadana, Yael Mathov, Yisroel Mirsky, Asaf Shabtai, Dominik Breitenbacher, and Yuval Elovici. N-BaloT:Networkbased Detection of IoT Botnet Attacks Using Deep Autoencoders. IEEE Pervasive computing, 17(3):12–22, 2018.
[Moore-2005]	Moore and Zuev, "Internet traffic classification using Bayesian analysis techniques", ACM SIGMETRICS, 2005
[Muehlstein-2017]	Muehlstein, J., Zion, Y., Bahumi, M., Kirshenboim, I., Dubin, R., Dvir, A. and Pele, "Analyzing HTTPS encrypted traffic to identify user's operating system, browser and application," In 14th IEEE Annu. Conf. Consum. Commun. Netw., pp. 1-6., 2017.
[Norbert-2021]	NYÁRI Norbert, "THE IMPACT OF QUANTUM COMPUTING ON IT SECURITY," Safety and Security Sciences Review, Vol 3, No 4, 2021.
[Ou-2019]	C. -W. Ou, F. -H. Hsu and C. -M. Lai, "Keep Rogue IoT Away: IoT Detector Based on Diversified TLS Negotiation," 2019 IEEE Intl Conf on Dependable, Autonomic and Secure Computing, Intl Conf on Pervasive Intelligence and Computing, Intl Conf on Cloud and Big Data Computing, Intl Conf on Cyber Science and Technology Congress (DASC/PiCom/CBDCom/CyberSciTech),

	Fukuoka, Japan, 2019, pp. 548-555, doi: 10.1109/DASC/PiCom/CBDCom/CyberSciTech.2019.00109.
[Palo Alto Networks]	<a href="https://www.paloaltonetworks.com/">https://www.paloaltonetworks.com/</a>
[Palo Alto Networks-1]	<a href="https://docs.paloaltonetworks.com/iot/iot-security-admin/detect-iot-device-vulnerabilities">https://docs.paloaltonetworks.com/iot/iot-security-admin/detect-iot-device-vulnerabilities</a>
[Palo Alto Networks-2]	<a href="https://docs.paloaltonetworks.com/iot/iot-security-admin/detect-iot-device-vulnerabilities/iot-device-vulnerability-detection">https://docs.paloaltonetworks.com/iot/iot-security-admin/detect-iot-device-vulnerabilities/iot-device-vulnerability-detection</a>
[Palo Alto Networks-3]	<a href="https://docs.paloaltonetworks.com/iot/iot-security-admin/discover-iot-devices-and-take-inventory/iot-device-discovery">https://docs.paloaltonetworks.com/iot/iot-security-admin/discover-iot-devices-and-take-inventory/iot-device-discovery</a>
[Palo Alto Networks-4]	<a href="https://www.paloaltonetworks.ca/resources/datasheets/zngbox">https://www.paloaltonetworks.ca/resources/datasheets/zngbox</a>
[Pan-2021]	J. Pan, "IoT Network Behavioral Fingerprint Inference with Limited Network Traces for Cyber Investigation," 2021 International Conference on Artificial Intelligence in Information and Communication (ICAIIC), Jeju Island, Korea (South), 2021, pp. 263-268, doi: 10.1109/ICAIIC51459.2021.9415273.
[Papadogiannaki-2021]	Eva Papadogiannaki and Sotiris Ioannidis. 2021. A Survey on Encrypted Network Traffic Analysis Applications, Techniques, and Countermeasures. ACM Comput. Surv. 54, 6, Article 123 (July 2022), 35 pages. <a href="https://doi.org/10.1145/3457904">https://doi.org/10.1145/3457904</a>
[Pashamokhtari-2021]	Pashamokhtari, N. Okui, Y. Miyake, M. Nakahara, H. H. Gharakheili, "Inferring Connected IoT Devices from IPFIX Records in Residential ISP Networks", IEEE LCN, Virtual Conference, Oct 2021.
[Pathmaperuma-2022]	M. H. Pathmaperuma et al, "CNN for User Activity Detection Using Encrypted In-App Mobile Data," Future Internet, vol. 14, (2), pp. 67, 2022.
[Quantum]	<a href="https://www.iotworldtoday.com/security/quantum-computing-offers-novel-security-for-iot-devices">https://www.iotworldtoday.com/security/quantum-computing-offers-novel-security-for-iot-devices</a>
[QUIC-Dataset-2018]	Q.: (2018), <a href="https://drive.google.com/drive/folders/1Pvev0hJ82usPh6dWDlz7Lv8L6h3JpWhE">https://drive.google.com/drive/folders/1Pvev0hJ82usPh6dWDlz7Lv8L6h3JpWhE</a>
[Qureshi-2021]	K.N. Qureshi, G. Jeon, F. Piccialli, Anomaly detection and trust authority in artificial intelligence and cloud computing, Comput. Netw. 184 (2021) 107647, <a href="http://dx.doi.org/10.1016/j.comnet.2020.107647">http://dx.doi.org/10.1016/j.comnet.2020.107647</a> , URL: <a href="https://www.sciencedirect.com/science/article/pii/S1389128620312664">https://www.sciencedirect.com/science/article/pii/S1389128620312664</a> .
[Rezaei-2019]	S. Rezaei and X. Liu, "Deep Learning for Encrypted Traffic Classification: An Overview," in IEEE Communications Magazine, vol. 57, no. 5, pp. 76-81, May 2019, doi: 10.1109/MCOM.2019.1800819.
[Rezaei-2020]	Shahbaz Rezaei and Xin Liu, "How to Achieve High Classification Accuracy with Just a Few Labels: A Semi-supervised Approach Using Sampled Packets," May 2020, <a href="https://arxiv.org/abs/1812.09761v2">https://arxiv.org/abs/1812.09761v2</a>
[Rhode & Schwarz Ipoque]	<a href="https://www.ipoque.com/">https://www.ipoque.com/</a>
[Rhode &	<a href="https://www.ipoque.com/blog/tackle-app-frenzy-with-dpi-driven-insights">https://www.ipoque.com/blog/tackle-app-frenzy-with-dpi-driven-insights</a>

Schwarz Ipoque-1]	
[Rhode & Schwarz Ipoque-2]	<a href="https://scdn.rohde-schwarz.com/ur/pws/dl_downloads/dl_common_library/dl_brochures_and_datasheets/pdf/1/Product-Flyer_vPACE_en_3683-7502-32_v0101.pdf">https://scdn.rohde-schwarz.com/ur/pws/dl_downloads/dl_common_library/dl_brochures_and_datasheets/pdf/1/Product-Flyer_vPACE_en_3683-7502-32_v0101.pdf</a>
[Rhode & Schwarz Ipoque-3]	<a href="https://www.ipoque.com/news-media/resources/ebooks/dpi-encrypted-traffic-visibility?cid=034">https://www.ipoque.com/news-media/resources/ebooks/dpi-encrypted-traffic-visibility?cid=034</a>
[Sandvine]	<a href="https://www.sandvine.com/">https://www.sandvine.com/</a>
[Sandvine-1]	<a href="https://www.sandvine.com/hubfs/Sandvine_Redesign_2019/Downloads/Whitepapers/sandvine_wp-encryption-and-dpi.pdf">https://www.sandvine.com/hubfs/Sandvine_Redesign_2019/Downloads/Whitepapers/sandvine_wp-encryption-and-dpi.pdf</a>
[Sandvine-2]	<a href="https://www.sandvine.com/hubfs/Sandvine_Redesign_2019/Downloads/2020/Whitepapers/Sandvine_WP_Advanced%20Traffic%20Classification.pdf">https://www.sandvine.com/hubfs/Sandvine_Redesign_2019/Downloads/2020/Whitepapers/Sandvine_WP_Advanced%20Traffic%20Classification.pdf</a>
[Sandvine-3]	<a href="https://www.sandvine.com/hubfs/Sandvine_Redesign_2019/Downloads/Whitepapers/Sandvine_WP_Encryption%20Use%20Cases%2020190625.pdf">https://www.sandvine.com/hubfs/Sandvine_Redesign_2019/Downloads/Whitepapers/Sandvine_WP_Encryption%20Use%20Cases%2020190625.pdf</a>
[Sarhan-2020]	M. Sarhan, S. Layeghy, N. Moustafa, M. Portmann, Netflow datasets for machine learning-based network intrusion detection systems, in: Big Data Technologies and Applications, Springer, 2020, pp. 117–135.
[Securolytics]	<a href="https://securolytics.io/">https://securolytics.io/</a>
[Securolytics-1]	<a href="https://www.cybersecurityintelligence.com/securolytics-7875.html">https://www.cybersecurityintelligence.com/securolytics-7875.html</a>
[Securolytics-2]	<a href="https://roi4cio.com/catalog/en/product/securolytics-iot-vulnerability-detection">https://roi4cio.com/catalog/en/product/securolytics-iot-vulnerability-detection</a>
[Shapira-2021]	Tal Shapira and Yuval Shavitt. Flowpic: A generic representation for encrypted traffic classification and applications identification. <i>IEEE Transactions on Network and Service Management</i> , 18(2):1218–1232, 2021.
[Shen-2021]	Meng Shen, Jinpeng Zhang, Liehuang Zhu, Ke Xu, and Xiaojiang Du. Accurate decentralized application identification via encrypted traffic analysis using graph neural networks. <i>IEEE Transactions on Information Forensics and Security</i> , 16:2367–2380, 2021.
[Shen-2023]	M. Shen et al., "Machine Learning-Powered Encrypted Network Traffic Analysis: A Comprehensive Survey," in <i>IEEE Communications Surveys &amp; Tutorials</i> , vol. 25, no. 1, pp. 791-824, Firstquarter 2023, doi: 10.1109/COMST.2022.3208196.
[Sivanathan-2019]	Sivanathan, H. H. Gharakheili, F. Loi, A. Radford, C. Wijenayake, A. Vishwanath and V. Sivaraman, "Classifying IoT Devices in Smart Environments Using Network Traffic Characteristics", <i>IEEE Transactions on Mobile Computing</i> , August 2019.
[Taylor-2016]	Vincent F Taylor, Riccardo Spolaor, Mauro Conti, and Ivan Martinovic. 2016. "Appscanner: Automatic fingerprinting of smartphone apps from encrypted network traffic," In 2016 IEEE European Symposium on Security and Privacy (EuroS&P). IEEE, 439–454.
[Taylor-2018]	Vincent F Taylor, Riccardo Spolaor, Mauro Conti, and Ivan Martinovic. 2018. "Robust smartphone app identification via encrypted network traffic analysis," <i>IEEE Transactions on</i>

	Information Forensics and Security 13, 1 (2018), 63–78.
[Towhid-2022]	M. S. Towhid and N. Shahriar, "Encrypted Network Traffic Classification using Self-supervised Learning," 2022 IEEE 8th International Conference on Network Softwarization (NetSoft), Milan, Italy, 2022, pp. 366-374, doi: 10.1109/NetSoft54395.2022.9844044.
[Tsogbaatar-2021]	E. Tsogbaatar, M.H. Bhuyan, Y. Taenaka, D. Fall, K. Gonchigsumlaa, E. Elmroth, Y. Kadobayashi, Del-IoT: A deep ensemble learning approach to uncover anomalies in IoT, Internet Things 14 (2021) 100391, <a href="http://dx.doi.org/10.1016/j.iot.2021.100391">http://dx.doi.org/10.1016/j.iot.2021.100391</a> , URL: <a href="https://www.sciencedirect.com/science/article/pii/S2542660521000354">https://www.sciencedirect.com/science/article/pii/S2542660521000354</a> .
[Ullah-2021]	Ullah, Imtiaz, and Qusay H. Mahmoud. "IoT Intrusion Detection Datasets." 2021, [Online]. Available: <a href="https://sites.google.com/view/iotdataset1">https://sites.google.com/view/iotdataset1</a> .
[Ullah-2022]	Ullah, Imtiaz, and Qusay H. Mahmoud. "An Anomaly Detection Model for IoT Networks based on Flow and Flag Features using a Feed-Forward Neural Network." 2022 IEEE 19th Annual Consumer Communications & Networking Conference (CCNC). IEEE, 2022.
[Wang-2015]	Q. Wang, A. Yahyavi, B. Kemme and W. He, "I know what you did on your smartphone: Inferring app usage over encrypted data traffic," 2015 IEEE Conference on Communications and Network Security (CNS), 2015, pp. 433-441, doi: 10.1109/CNS.2015.7346855.
[Wang-2018]	P. Wang, F. Ye, X. Chen and Y. Qian, "Datanet: Deep Learning-Based Encrypted Network Traffic Classification in SDN Home Gateway," in IEEE Access, vol. 6, pp. 55380-55391, 2018.
[Wu-2022]	Y. Wu, H. -N. Dai and H. Tang, "Graph Neural Networks for Anomaly Detection in Industrial Internet of Things," in IEEE Internet of Things Journal, vol. 9, no. 12, pp. 9214-9231, 15 June 15, 2022, doi: 10.1109/JIOT.2021.3094295.
[Zahid-2022]	Zahid, Hafiz Muhammad, et al. "A Framework for Identification and Classification of IoT Devices for Security Analysis in Heterogeneous Network." Wireless Communications and Mobile Computing 2022 (2022).
[Zhang-2011]	Fan Zhang, Wenbo He, Xue Liu, and Patrick G. Bridges. 2011. Inferring users' online activities through traffic analysis. In Proceedings of the fourth ACM conference on Wireless network security (WiSec '11). Association for Computing Machinery, New York, NY, USA, 59–70. <a href="https://doi-org.proxy.library.carleton.ca/10.1145/1998412.1998425">https://doi-org.proxy.library.carleton.ca/10.1145/1998412.1998425</a>
[Zhao-2019]	Z. Zhao, S. Cerf, R. Birke, B. Robu, S. Bouchenak, S. Ben Mokhtar, L.Y. Chen, Robust anomaly detection on unreliable data, in: 2019 49th Annual IEEE/IFIP International Conference on Dependable Systems and Networks, DSN, 2019, pp. 630–637, <a href="http://dx.doi.org/10.1109/DSN.2019.00068">http://dx.doi.org/10.1109/DSN.2019.00068</a> .
[Vishwakarma-2022]	M. Vishwakarma and N. Kesswani, "DIDS: A Deep Neural Network based real-time Intrusion detection system for IoT," Decision Analytics Journal, vol. 5, p.100142, 2022.

## Acronym/Glossary

Acronym	Meaning
ARP	Address Resolution Protocol
BoT	A software application that runs automated tasks (scripts) over the Internet, usually with the intent to imitate human activity on the Internet, such as messaging, on a large scale. [ <a href="https://en.wikipedia.org/wiki/Internet_bot">https://en.wikipedia.org/wiki/Internet_bot</a> ]
C4.5	C4.5 is an algorithm used to generate a decision tree developed by Ross Quinlan. [ <a href="https://en.wikipedia.org/wiki/C4.5_algorithm">https://en.wikipedia.org/wiki/C4.5_algorithm</a> ]
CAGR	Compounded Annual Growth Rate
CNN	Convolutional Neural Network
CIC	Canadian Institute for Cybersecurity
CSE	Communications Security Establishment
DAE-EPNN	Deep AutoEncoder-EnsembleProbabilistic Neural Networks
DeL-IoT	Deep ensemble Learning model-based framework for IoT anomaly detection
DIDS	Deep Neural Network-based Intrusion Detection System
DFT	Discrete Fourier Transform
DNS	Domain Name System
DPI	Deep Packet Inspection
DR	Detection Rate
desIP	Destination IP
desPort	Destination Port
DoS	Denial of Service
ECE	ECN-Echo, used to echo back the congestion indication
ECN	Explicit Congestion Notification
ENTA	Encrypted Network Traffic Analysis
ESNI	Encrypted Server Name Indication
ETI	Encrypted Traffic Intelligence
$F_1$ Score	It is the harmonic mean of a system's precision and recall values, calculated as follows: $2 \frac{Precision \times Recall}{Precision + Recall}$
FCBF	Fast Correlation-Based Filter
FDR	False Detection Rate
FIN	A message that triggers a graceful connection termination between a client and a server [ <a href="https://www.baeldung.com/cs/tcp-fin-vs-rst#:~:text=FIN%3A%20a%20message%20that%20triggers,a%20client%20and%20a%20server">https://www.baeldung.com/cs/tcp-fin-vs-rst#:~:text=FIN%3A%20a%20message%20that%20triggers,a%20client%20and%20a%20server</a> ].
FFNN	Feed Forward Neural Network
FPR	False Positive Rate
GP	Genetic Programming
HDNN	Hierarchical Deep Neural Network
HTTPS	Hypertext Transfer Protocol Secure

ICMP	Internet Control Message Protocol
IIoT	Industrial IoT
iOS	iPhone Operating System
IoT	Internet of Things
IP	Internet Protocol
IPFIX	IP Flow Information eXport
ISP	Internet Service Provider
IT	Information Technology
KDA	Kernel Discriminant Analysis
LDA	Linear Discriminant Analysis
LEA	Law Enforcement Agency
LSTM	Long Short-Term Memory
MAC	Media Access Control
MCC	<p>Matthews Correlation Coefficient</p> <p>It measures the quality of a classification, showing the correlation agreement between the observed values and the predicted values. Its equation is as follows:</p> $\frac{TP \times TN - FP \times FN}{\sqrt{(TP+FP) \times (TP\_FN) \times (TN+FP) \times (TN\_FN)}}, \text{ where}$ <ul style="list-style-type: none"> <li>• TP: True Positive</li> <li>• TN: True Negative</li> <li>• FP: False Positive</li> <li>• FN: False Negative</li> </ul>
MQTT	Message Queuing Telemetry Transport, an OASIS standard messaging protocol for the Internet of Things (IoT) [ <a href="https://mqtt.org/">https://mqtt.org/</a> ].
ML/DL	Machine Learning/Deep Learning
MUD	Manufacturer Usage Description
N-BaloT	Network-Based Approach for IoT
NAT	Network Address Translation
NF	Netflow
NTP	Network Timing Protocol
OT	Operational Technology
P2P	Peer-to-Peer
PCA	Principal Component Analysis
PCAP	Packet Capture
PNN	Probabilistic Neural Network
PSH	Push flags for instructing the operating system to send or receive data immediately, respectively at the source or receiver
QoE	Quality of Experience
QUIC	QUIC is a new multiplexed transport built on top of UDP. HTTP/3 is designed to take advantage of QUIC's features, including lack of Head-Of-Line blocking between streams [ <a href="https://www.chromium.org/quic/">https://www.chromium.org/quic/</a> ].
RF	Random Forest

RFC	Request for Comments
RFE	Recursive Feature Elimination
RRSE	Root Relative Square Error
RSA	RivestShamirAdleman is a public-key cryptosystem that is widely used for secure data transmission. [ <a href="https://en.wikipedia.org/wiki/RSA_(cryptosystem)">https://en.wikipedia.org/wiki/RSA_(cryptosystem)</a> ]
RST	A message that aborts the connection (forceful termination) between a client and a server [ <a href="https://www.baeldung.com/cs/tcp-fin-vs-rst#:~:text=FIN%3A%20a%20message%20that%20triggers,a%20client%20and%20a%20server">https://www.baeldung.com/cs/tcp-fin-vs-rst#:~:text=FIN%3A%20a%20message%20that%20triggers,a%20client%20and%20a%20server</a> ].
SAE-EPNN	Stacked AutoEncoder-Ensemble Probabilistic Neural Networks
SBB-GP	Symbiotic Bid-based GP
SCP	Secure Copy
ScrlP	Source IP
ScrPort	Source Port
SDN	Software-Defined Network
SFTP	Secure File Transfer Protocol
SOC	Security Operations Center
SotA	State-of-the-Art
SSH	Secure Shell
SSL	Secure Sockets Layer
SVM	Support Vector Machine
TCP	Transport Control Protocol
TLS	Transmission Layer Security
TON-IoT	IoT/IoT datasets collected from Telemetry data, Operating systems data and Network data.
tos	Type of Service
TRC	Traffic Rate Change
Trust Metric	<p>From [Pashamokhtari-2021], the intuition behind the trust metric is to check the number of expected flows from a device with the number of flows that are indeed classified as that of the device type during the monitoring period.</p> <p>The raw measure of trust is computed as: <math>T_L = \frac{N_{o,L}}{N_{e,L} \times D_{e,L}}</math>,</p> <p>where <math>N_{o,L}</math> is the number of observed flows predicted as class <math>L</math>; while <math>N_{e,L}</math> is expected number of record of class <math>L</math> and <math>D_{e,L}</math> is the expected rate of discarded records during training.</p> <p>The normalized trust is defined as: <math>T_{norm,L} = \exp\left(-\frac{(T_L - 1)^2}{2\sigma_{T_L}^2}\right)</math>,</p> <p>where <math>\sigma_{T_L}^2</math> is the standard-deviation of <math>T_L</math> computed during the training phase.</p>
UDP	User Datagram Protocol
VPN	Virtual Private Network