

VĚDECKÉ SPISY VYSOKÉHO UČENÍ TECHNICKÉHO V BRNĚ

*Edice Habilitační a inaugurační spisy, sv. 757*

*ISSN 1213-418X*

**Václav Oujezský**

**SELECTED TECHNIQUES INVOLVED  
IN NETWORK TRAFFIC ANALYSIS**

**BRNO UNIVERSITY OF TECHNOLOGY**  
**Faculty of Electrical Engineering and Communication**  
**Department of Telecommunications**

**Ing. Václav Oujezský, Ph.D.**

**SELECTED TECHNIQUES INVOLVED  
IN NETWORK TRAFFIC ANALYSIS**

**VYBRANÉ TECHNIKY ANALÝZY SÍŤOVÉHO PROVOZU**

**TEZE HABILITAČNÍ PRÁCE  
V OBORU TELEINFORMATIKA**



**BRNO 2023**

**KEYWORDS**

Analysis, Data, Network, Traffic, Techniques

**KLÍČOVÁ SLOVA**

Analýza, Data, Provoz, Síť, Techniky

**MÍSTO ULOŽENÍ PRÁCE**

Vědecké oddělení, Faculty of Electrical Engineering and Communication, Brno  
University of Technology, Technická 3058/10, 616 00 Brno

© Václav Oujezský, 2023

ISBN 978-80-214-6163-5

ISSN 1213-418X

# CONTENTS

INTRODUCTION . . . . .	5
1 THESIS OVERVIEW . . . . .	6
1.1 Goals . . . . .	6
1.2 Contribution and Relation to Author’s Publications . . . . .	6
1.3 The Organisation of the Thesis . . . . .	7
2 THE STATE OF THE ART IN NETWORK TRAFFIC ANALYSIS . . . . .	7
3 DATA SOURCES FOR TRAFFIC ANALYSIS . . . . .	11
3.1 Network Traffic Flows . . . . .	11
3.1.1 <i>Network Traffic Flows and Cloud Environment</i> . . . . .	12
3.2 Full Data Packets . . . . .	12
3.2.1 <i>Full Packet Data and Cloud Network</i> . . . . .	13
3.3 Specific Data Sources . . . . .	13
3.3.1 <i>Sensors, Wireless, and Internet of Things Networks</i> . . . . .	13
3.3.2 <i>Dedicated Hardware</i> . . . . .	14
3.4 Datasets . . . . .	14
4 DATA PREPARATION AND PROCESS STREAMLINING TECHNIQUES . . . . .	15
4.1 Data Preparation . . . . .	15
4.2 Data Classification and Clustering Techniques . . . . .	16
4.2.1 <i>The Role of Evolutionary Algorithms</i> . . . . .	18
4.2.2 <i>Process Streamlining by Using Genetic Algorithms</i> . . . . .	19
5 SELECTED TRAFFIC ANALYSIS TECHNIQUES . . . . .	20
5.1 Location Accuracy Analysis . . . . .	21
5.2 Network Traffic Behavior Analysis . . . . .	22
5.3 Machine Learning Analysis . . . . .	23
5.4 Analysis using Artificial Immune Systems . . . . .	25
5.4.1 <i>Mapping selected problem into the AIS</i> . . . . .	27
6 NETWORK DATA ANALYSIS AND REPORTING . . . . .	28
7 VERIFICATION AND INTERPRETATION OF RESULTS . . . . .	31
7.1 General approach . . . . .	31
CONCLUSION . . . . .	32

## CHARACTERISTICS OF THE AUTHOR



**Vaclav Oujezsky** was born in Brno, The Czech Republic. He currently works as an assistant professor at Masaryk University and a researcher at the Brno University of Technology, Department of Telecommunications. He is working actively on projects of security and transport networks. His research interests include computer networks, network programming, and software-defined networking. He focuses on network behavior, intelligent networks, network analysis, and security.

**PROFESSIONAL EXPERIENCE:** Vaclav Oujezsky studied from 2008–2013 for a Master’s degree in teleinformatics at the Brno University of Technology. In the years 2013–2017, he received a Ph.D. in the field of teleinformatics. From 2006 to 2014, he worked as a Network Engineer at T-Systems Czech Republic, from 2014 to 2016 as a Senior Network Engineer at T-Mobile CZ, from 2016 to 2021 as a Senior Network Engineer at IBM Client Innovation Center Brno, and at the same time as a Researcher at the Department of Telecommunications, Brno University of Technology, and from 2021 as an academic staff member at the Masaryk University. Over the years, he has obtained several certificates, including CCNA Cisco Certified Network Associate, CCNP Cisco Certified Network Professional, IBM Mentor IBM Profession – Technical Specialist, NSA Certificate of Security Clearance – Level: Confidential, and others. He has participated in five grant projects and is a co-author of a number of scientific articles and conference papers. He is a leading expert in the field of modern information systems.

**PROFESSIONAL FOCUS:** The author’s research activities are focused mainly on analyzing data traffic in communication networks. He has participated in or is participating in five scientific projects of the Ministry of the Interior of the Czech Republic (MVCR), of which two projects as the principal investigator and one project of the Technology Agency of the Czech Republic (TA CR) as a co-principal investigator. He is a member of the working group of the European COST Action. He participated in the implementation of three functional samples and six software results. He is actively involved as a reviewer of research projects for MVCR or TACR and is a member of the scientific conference committee and workshop co-organizer. The author has published research results in a total of 51 articles or scientific conferences, of which 41 papers are indexed in WoS, 40 are indexed in Scopus, and 12 are in impacted journals. His articles have received 99 citations according to WoS and 132 citations according to Scopus, and his H-Index is 6.

## INTRODUCTION

In recent years, the study of large-scale networks has become widespread. Especially with the Internet boom, the first opportunity to study large-scale networks has emerged. Today's wide area networks are increasingly a part of society, and society is more dependent on them than ever before – their functioning results from a combination of many technologies and many principles. In the search for efficiencies in their combined functioning and the search for ways to provide protection, researchers have begun to systematically analyze and characterize the dynamic patterns and evolution of these complex and highly heterogeneous structures. Their properties are usually manifested in connectivity patterns characterized by large fluctuations, scale-free properties, and non-trivial correlations such as high clustering and hierarchical ordering. The large size and dynamic structure of complex networks are closely related to graph theory and the ability to characterize the dynamic evolution of the basic elements of such a system. Advances in research on complex networks have sparked interest in the possible implications and consequences of their operation for the most important questions concerning various dynamic processes and functions.

The issue of computer network security is not fading in importance; on the contrary, interest in addressing security issues in modern communication networks is growing. Unsolicited traffic nowadays significantly affects the security of society. The current trend is to use a hybrid approach, where attackers typically mix several operations together to create many attack vectors. During these operations, attack vectors and protocol signatures are altered to deceive automated mitigation devices. Increasingly, mobile clients and smart devices are at the forefront of attackers' minds. The need to develop security devices that can respond in real-time to attacks and network traffic anomalies is undoubtedly a priority.

On the other hand, providing a general account of the study of traffic networks is impossible since each research domain is very specific. The research I have conducted after I obtained my Ph.D. degree in Telecommunications has mainly had in view two areas that involve traffic analysis in passive optical networks and methods used for traffic behavior analysis. The research I have done is related to projects in which I have personally participated or proposed and I am a co-investigator; details of some of the projects are published in [1]. The research itself focused on two fundamental questions. The design of a system to enable the detection and analysis of traffic in communication networks is partially and recently focused on passive gigabit networks. The second question concerns the possibility of analyzing and verifying the data obtained from such network traffic. As a result of the research conducted, the habilitation thesis entitled „*Selected techniques involved in network traffic analysis*“ unfolds in the following directions: research pertaining to methods for network traffic analysis in gigabit passive optical networks and the use of sub-parts of artificial intelligence, such as genetic algorithms, machine learning or data processing. One part of the research also includes the mentioned system design for data capture and analysis. A specific area of research related to passive optical networks presents complex networks different from other types of networks based on Ethernet

technology. The popularity and use of these networks are constantly growing and, therefore, need to be addressed from all aspects. The fundamental shortcomings of traffic analysis in optical networks were presented in [2]. Several key aspects that are directly related to the security of passive optical networks are mentioned.

My research work has resulted in algorithms and methods applicable to traffic analysis in several types of networks. Another result is a concept of a comprehensive model of a system for traffic analysis. Both results of which are the content of this habilitation thesis. None of the results or proposed solutions presented in this thesis were published in my Ph.D. thesis or any previous theses. It contains, of necessity, only some references or insights that support the continuous progress of scientific work in the related area.

## **1 THESIS OVERVIEW**

This chapter contains an overview of this thesis. Section 1.1 provides a summary of its primary objectives. In Section 1.2, the author's publishing activities and contribution to the scientific subject of the thesis are briefly discussed. Finally, Section 1.3 introduces the organizational structure of the thesis.

### **1.1 GOALS**

In this habilitation thesis, the main topic is traffic and data analysis methods in complex and heterogeneous networks and techniques that provide and support data analysis. The main idea of the work is laid out in several parts. This reflects the distinct focus of each of the sections that follow and identifies the relevant in-depth publications. I also provide a quick summary of the most pertinent connected works in each area here, and details are then given in the particular included publications. The aim of the thesis is to present the latest new knowledge in the field of network analysis techniques obtained from other experts and especially from the author's research activities.

### **1.2 CONTRIBUTION AND RELATION TO AUTHOR'S PUBLICATIONS**

The thesis content is written considering both pedagogical and scientific contributions. In order to provide a foundational understanding, an overview of the state of the art, and familiarity with the most recent scientific results, it should be helpful not only to specialists in the field of data traffic analysis but also to students interested in this topic. The scientific contributions of the author to the topic are listed at the beginning of each chapter for better transparency.

The present habilitation thesis manuscript covers the scientific progress I have achieved in the years 2017 to 2022, by exploring different topics in traffic analysis, techniques, and methods of detection of malicious traffic or detection of abnormal traffic, especially in passive optical networks, but not limited to this type of network. The main results and proposals presented in the thesis have been published in various international journals with impact factors, such

as Sensors *International Standard Serial Number* (ISSN): 1424-8220, Electronics ISSN: 2079-9292, or are indexed by Web of Science and Scopus. In addition, I am also the primary author or coauthor of various publications of other topics related to telecommunication networks.

The research results described in the thesis are partially the subject of recent research projects *Ministry of the Interior of the Czech Republic* (MV ČR) 2017-2019 – “Detection of security threats on the active components of critical infrastructures” (setting up, co-investigator, researcher, programmer), MV ČR 2017-2020 – “Reduction of security threats at optical networks” (setting up, co-investigator, researcher, programmer), MV ČR 2019-2022 – “Deep hardware detection of network traffic of next-generation passive optical network in critical infrastructures” (setting up, co-investigator, researcher, programmer), *Technology Agency of the Czech Republic* (TA ČR) 2019-2024 – “Decentralized Control of Distribution System” (research team member), and from the activities related to the foreign collaboration, as is COST Action CA20120 2021-2025 (working group member) and last but not least, from cooperation with Assoc. Prof. Vladislav Škorpil, Dr. Tomáš Horváth, Brno University of Technology, or from cooperation with industrial partners. Current projects in the role of the lead investigator are MV ČR VK01030030 2023-2026 – “Data backup and storage system with integrated active protection against cyber threats”, and MV ČR VK01030152 2023-2024 – “Android federated learning framework for emergency management applications”.

### 1.3 THE ORGANISATION OF THE THESIS

The thesis is divided into seven basic chapters. The thesis is organised as follows. In Chapter 1, the thesis overview is given, including the goals and contributions in related research, projects, and publications. Chapter 2 presents the state-of-the-art in the field of network traffic analysis, giving a general overview of the techniques used for network analysis. Chapter 3 presents the topic related to the source of data for network traffic analysis and the data sources currently used for network traffic analysis and processing techniques. Chapter 4 discusses the data preparation and selected process streamlining techniques and algorithms used. Chapter 5 deals with data traffic analysis techniques such as network traffic behavior analysis, deep packet inspection, network frame structure analysis, or the use of artificial intelligence in traffic analysis. Chapter 6 summarizes current data monitoring concerns and discusses the possible monitoring of passive optical networks. Chapter 7 proposes a template for the verification and interpretation of traffic analysis results. The conclusion provides a review of the thesis and planned short and long-term future works.

## 2 THE STATE OF THE ART IN NETWORK TRAFFIC ANALYSIS

In computer science, the analysis can be object-oriented or syntactic. In networking, *Network Traffic Analysis* (NTA) refers to the discovery and understanding of events that occur



in the operation of network elements for the purpose of operation and protection. From the latest research and experience, it can be concluded that network analysis generally, from a high-level point of view, consists of behavioral analysis or appearance analysis. These two can be analyses of an individual element or neighborhood (the surrounding environment or elements) and its connections (network traffic flows, network traffic connections), and these two again in real or relative time and space. The traffic analysis is performed for the purpose of finding an anomaly, to detect unsolicited traffic, or to detect malicious code content, etc (an analysis is, therefore, subject to subsequent detection). Therefore, anomaly detection is the process by which anomalies are detected against the relevant data. In contrast, an abnormality is a behavioral dysfunction, that is, a change in its behavioral characteristics that causes deviations. Although academics are looking for other uses for the data in traffic flows, the fundamental purpose of these data is to monitor the network's performance.

First, the data must be obtained in some way in the form of some datasets or other records, such as data flow protocols, and then processed into a suitable form for further processing. Second, the data must be prepared for the analysis itself, and techniques for data preprocessing or re-sorting must be applied. The actual data analysis or extraction follows this step. The last step has to include the evaluation of the analysis results, which acts as feedback for the analysis algorithms and processes themselves. The basic data analysis process is shown in Figure 2.1. Currently, there are many approaches to the analysis itself. Among the general approaches discussed above, the following types are defined, but not limited to:

- **Host-based analysis** – refers to the collection and examination of information from a host, including but not limited to live memory collection and analysis, traffic analysis, file carving and recovery, and data analysis.
- **Network-level-based analysis** – refers to the process of examining a network's availability and activity. Tracking the data flow across various network segments includes determining what data are being sent when, and how. It can also be performed by:
  - **Flow-based traffic analysis** – in the sense of using any method to examine traffic flows (streams) between a source and a destination. These techniques frequently focus on looking at network flows, typically described as a series of packets sent over a certain amount of time and organized by protocol, source address, source port, destination address, and destination port.
  - \* **Statistical analysis** – for traffic analysis, either multivariate models or mod-

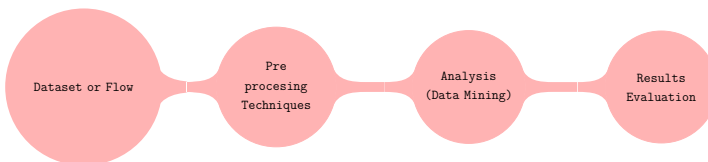


Fig. 2.1: The Generic Structure of Network Traffic Data Analysis.

els based on available statistics such as data entropy, compression, or measurements of the mean and variance of predefined profiles are used. Cluster analysis is also often used in statistical methods. In this method, patterns of normal network traffic are also used. The basic principles of those are simple average, data entropy measurement, and similarity comparison.

- \* **Graph-theory-based analysis** – in the sense of using applied graph theory for analysis to calculate network measures.
- \* **Signature-based analysis** – typically involves payload or deep packet inspection techniques and looking for specific patterns. Patterns can be extracted manually or automatically using an automatic signature extraction mechanism.

There are now three main research areas, as noted in [3], and hence three traffic inspection techniques. The first is the *Deep Packet Inspection* (DPI) approach, which separately decrypts and inspects each packet. Behavior analysis is the second option. This method uses several techniques. For example, it counts the number of packets transmitted or measures flow parameters such as the interval between packets, etc., for the purpose of discovering a pattern of behavior without needing to know the content of the data. In addition to the methods mentioned, the fingerprinting technique is also an option. This technique takes advantage of data that can be seen during the early stages of an encrypted connection.

Other techniques beyond these three can include various methods to gather more information in combination to make traffic analysis effective. Such as the use of *Simple Network Management Protocol* (SNMP) and device statistics, or other new approaches such as *Data Model-Driven Management* (DMDM) using the *Yet Another Next Generation* (YANG) modeling language for the network configuration protocol, defined by *Request for Comments* (RFC) 7950 [5]. This can be part of network telemetry, which is also an option to gather useful information for network analysis. Although the SNMP protocol has historically been very effective at monitoring, it has some drawbacks. Telemetry provides an alternative way of addressing many of the shortcomings of older monitoring tools. The topic of network telemetry is best described by the Network Telemetry Framework RFC 9232 [4]. In a generic term, in the networking world, telemetry is an automated communication process by which data and

	Management Plane	Control Plane	Forwarding Plane
Data Configuration and Subscribe	gNMI, NETCONF, RESTCONF, SNMP, YANG-Push	gNMI, NETCONF, RESTCONF, YANG-Push	NETCONF, RESTCONF, YANG-Push
Data Generation and Process	MIB, YANG	YANG	IOAM, PSAMP, PBT, AM
Data Encoding and Export	gRPC, HTTP, TCP	BMP, TCP	IPFIX, UDP

Fig. 2.2: The work mapping of network telemetry, as is in RFC 9232 [4].

measurements are collected at remote devices and transmitted to a monitoring device. Today, many types of telemetry implementations and protocols are used, for example, the *In-band Network Telemetry* (INT) dataplane specification [6] or *In Situ Operations, Administration, and Maintenance* (IOAM) defined by RFC 9197 [7]. The purpose of RFC 9232 is to provide an entire telemetry framework, shown in Figure 2.2.

However, the topic of network automation and network telemetry is beyond the scope of this thesis. Basically, INT may be used for monitoring from the hardware level by data plane programming for service quality monitoring or microburst detection for latency demand application<sup>1</sup>. Therefore, the data analysis options can be summarized in the following techniques.

- Techniques using **traffic flow protocols**, discussed in Section 3.1.
- Techniques using **deep packet inspection** and **full packet capture** (PCAP), discussed in Sections 3.2.
- Techniques using **statistic collections** or specific data collections, discussed in Section 3.3.
- Techniques using a **mix of approaches**.

There are many approaches to obtaining traffic characteristics, in other words, a kind of basic traffic analysis. The basic characteristic is related to the performance of the network layer. Characteristics such as delay, throughput, or packet loss are measured and analyzed. One of the simplest ways is to calculate a moving average, which can be applied, for example, to flow analysis or to delay analysis, respectively. In the case of *Simple Moving Average* (SMA), it is an unweighted average of  $n$  numbers in a time series. There are also other variants based on the SMA. These are cumulative average, weighted moving average, exponential moving average, and other modified ones. The *Exponential Moving Average* (EMA) is used in the *Round-Trip Time* (RTT) calculation of bidirectional network traffic delay by the PING (*Packet InterNet Groper*) program. The use of only a moving average is not sufficient enough to reflect the various aspects of traffic analysis. For example, to measure the “burst” level of traffic streams. Entropy in information technology is also referred to as *Shannon Entropy* ( $H$ ) after the author Claude Elwood Shannon. It is a calculation of the quantity of information for some whole phenomenon. If this phenomenon is an occurrence of a signal element, then its unit is [Sh/element].

Many traffic analysis techniques are based on the concept of distance. The most well-known metric is the Euclidean distance. As such, it is widely used to measure continuity or similarity. In the field of traffic analysis, algorithms based on biologically inspired methods, commonly known as *Artificial Intelligence* (AI), are also used. These include, in particular, neural networks, computations using evolutionary *Artificial Immune System* (AIS) algorithms, and others. For example, evolutionary algorithms have been used to improve existing traffic analysis methods. The authors in [8] used *Genetic Algorithm* (GA) for the extraction of network traffic components. The author in [9] used the Genetic Algorithm to calculate the Euclidean distance

---

<sup>1</sup>GÉANT: <https://wiki.geant.org/display/NETDEV/DPP>

matrices. AIS was used by the author in [10] to detect anomalous traffic and by the author in [11] to analyze *Gigabit Passive Optical Network* (GPON) frames.

The new paradigms for the use of AI are not only for their use in traffic analysis, but also for the creation of a new communication model concerning the use of a semantic learning model to specify the values of the communication channel applied to the Shannon theorem [12].

Network traffic analysis is also a matter of hardware, which goes hand in hand with the increasing speed of computer networks. Today we are already commonly facing a network speed of 100 Gb/s. Commercial equipment can be advantageously used, such as a high-speed analyzer capable of storing large amounts of data or high-speed traffic generators. In addition, specialized hardware elements such as *Field-programmable Gate Array* (FPGA) cards capable of handling large amounts of data, or devices equipped with *Graphics Processing Units* (GPU)s or *Tensor Processing Unit* (TPU)<sup>2</sup>s, or cloud environments, can be used in combination for the analysis itself or the development of new techniques.

### 3 DATA SOURCES FOR TRAFFIC ANALYSIS

This chapter is composed of the author's own results and supplemented with theoretical knowledge on the subject from other sources. It deals with the problem of data acquisition for traffic analysis. The first Section 3.1 discusses network flows and their versions and outlines the topic of cloud solutions. Section 3.2 provides an overview of the techniques used to extract full traffic data from network traffic. Specific data sources and the author's published results are discussed in Section 3.3, and data sets as a data source are discussed in Section 3.4.

#### 3.1 NETWORK TRAFFIC FLOWS

**The authors' contribution:** Software: <sup>3</sup> Network Analyzer – detects and analyzes incoming NetFlow messages (versions 1, 5, and 9 in the latest version) of network devices that support them. It works in *Command Line Interface* (CLI). The output file is a database of information and analysis of the overall UNIX time duration of each reported traffic. Software has been developed to work with Python version 3 and greater, designed for the Windows operating system. Another contribution in this area is the leading of student theses, for example, on: Design and implementation of the network collector<sup>4</sup>.

---

The most commonly used data format for traffic analysis is, undoubtedly, network monitoring traffic flows. Network monitoring traffic flows are used to monitor performance, security, and other factors. Network flows differ significantly from packet captures in that they provide only details about the transmission event, such as source and destination addresses, the

---

<sup>2</sup>The Tensor Processing Unit is a specialized processing unit (hardware) developed by Google LLC primarily for their project as is, for example AlphaGo by providing better machine learning results and they released it into the cloud so other people can use it in their machine learning projects.

<sup>3</sup>Network Research Group – <https://nsr.utko.feec.vutbr.cz/software.php>

<sup>4</sup>Jaroslav Bořela: <http://bit.ly/30kGc8i>

network protocol used, the quantity of data transported, and other details, rather than information on the content of sent (user) data (payload).

The widely adopted NetFlow network flow monitoring standard was invented by Cisco Systems, Inc. Several multiple versions of NetFlow are now in use, most notably NetFlow 5, NetFlow 9, defined in RFC 7011 [13]. The network flow protocol in some modifications is used by a number of other vendors, such as Juniper (Jflow); 3Com/HP, Dell, and Netgear (s-flow); Alcatel-Lucent (Cflow); Ericsson (Rflow). Another protocol is *IP Flow Information Export* (IPFIX), which is an *Internet Engineering Task Force* (IETF) protocol, defined in RFC 5101 [14], later extended in RFC 7011 [13].

### 3.1.1 Network Traffic Flows and Cloud Environment

Cloud networks are different in the sense of virtualization, and the situation for flow analysis is not as straightforward here. Cloud systems are made up of a complex set of different technologies. A datacenter can be thought of as a multilayered system, similar to the conceptual *Open Systems Interconnection* (OSI) model. Each of these layers has specific ways to extract data for further analysis. Typically, a cloud service provider has the authority and power of attorney over Layer 1 to Layer 3 devices. From the following upper Layer 3 to Layer 6 it is then a combination where funds are provided to individual lessees, customers. Customers can then build a virtualized network as an overlay layer on the physical layer. Or they can run virtualized operating systems over which they have management. But customers have no visibility into the provider-managed layer. So those are the two basic views of the data source. So in cloud services, data must be obtained at higher layers in different ways than in legacy networks. The data come from the overlay layer if it is data from the customer's perspective. The situation around flows is different, and specialized tools such as CloudWatch flow log AWS Lambda or AZURE FlowLogs operating at network layer 4 have been developed. These techniques have some limitations in terms of updating flows by seconds or in their total amount. The other option is to use JSON blobs and visibility data (similar to SPLUNK) generated by intelligent applications.

## 3.2 FULL DATA PACKETS

Whole data is needed when deeper analysis of traffic patterns is performed, such as DPI, data pattern matching, or signature search. The process of obtaining the data is called packet acquisition. The most well-known is the PCAP format. It is the acronym for packet capture API. The Windows *Operating System* (OS) implementation is known as WinPcap, while the Unix OS implementation is known as libpcap. Today, there are manufacturers of FPGA network cards that support the libcap library and can be used in high-speed networks. There are wrappers for various programming languages available, as is the python-libpcap.

The next option is to use network *Test Access Point* (TAP) or *Switch Port Analyzer* (SPAN). The SPAN is a specific tool included in a network switch that copies Ethernet frames passing through switch ports and sends these frames out to a specific port to the network analyzer.

### 3.2.1 Full Packet Data and Cloud Network

As in legacy network, a virtual network TAP can be used. Azure virtual network TAP allows to continuously stream virtual machine network traffic to a network packet collector or analytic tool. The other possible way is to use cloning of network traffic by using specialized network devices as is traffic manager. *Encapsulated Remote Switch Port Analyzer* (ERSPAN) is also another way how to collect traffic from virtualized routers to network probes in an AWS cloud based network.

### 3.3 SPECIFIC DATA SOURCES

**The authors' contribution:** Design and implementation of a system for analysis of GPON and 10 Gigabit-capable PON (XGPON) frames from downstream and upstream traffic. The design includes parsing and processing of the bit stream obtained from the communication between *Optical Line Terminator* (OLT) and *Optical Network Unit* (ONU) using a splitter and an FPGA network card. The result is a custom data parser software that parses the data and sends it in *JavaScript Object Notation* (JSON) format to Apache Kafka for further processing. In addition, hosting a *Réseaux IP Européens* (RIPE) probe<sup>5</sup> to investigate the possibilities of anomaly detection using the network probe and developing an application as part of the thesis supervision.

---

Another way to obtain data and information about network traffic besides using full packet capture or using network flows is to use proprietary approaches and data formats. One representative is the JSON used by the RIPE community. The Resource Request *Application Programming Interface* (API) is used to submit requests for Internet number resources to the RIPE NCC and the supported requests at the time being are *Autonomous System* (AS) Number Assignment, IPv4 First Allocation, IPv6 First Allocation, or IPv6 *Provider Independent* (PI) Assignment. The information is retrieved from the RIPE databases and sent back to the user, who can then deserialize it and work with it as part of an analysis application to detect anomalies in networks. The traffic analysis can be trace-route based, such as route changes in neighboring connections, delay in entire neighboring network, RTT increase, or ping based, such as anchor<sup>6</sup> delay per country or anchor down, etc.

#### 3.3.1 Sensors, Wireless, and Internet of Things Networks

**The authors' contribution:** Research in IoT security and securing end-device communications and cloud applications[15].

---

A special case, not necessarily separate from classical networks, are sensors, specifically *Wireless Sensor Networks* (WSN) and *Internet of Things* (IoT) networks. When determining

---

<sup>5</sup><https://www.fi.muni.cz/~oujezsky/probe.html>

<sup>6</sup><https://atlas.ripe.net/about/anchors/>

performance indicators and making subsequent decisions to enhance network performance, modeling the behavior of the networks becomes necessary. This topic is very extensive, so in this thesis, it is limited to the overview of data acquisition. The possibilities of acquiring data directly from sensors and the principle of transmission are more related to information acquisition than data acquisition for the purpose of traffic analysis as such. In terms of analysis of the network traffic itself, parameters such as path loss, delay, throughput or sensor network life time can be monitored. Basic mathematics for IoT system network is tied with the rules of classical networks.

### 3.3.2 Dedicated Hardware

**The authors' contribution:** a project focused on the development of a programmable FPGA network card for gigabit-capable passive optical networks. The network card is constructed to analyze GPON frames and check the correctness of communication between an optical line terminator and an optical network unit directly in the optical domain. The GPON networks use a specific encapsulation method for Ethernet frames and control messages. Because of the network card operates directly in the optical domain, it is possible to provide real-time analysis of headers of GPON transmission convergence layer [16].

---

As mentioned in the introductory chapter, analytical tools can be supplemented with specific programmable hardware that then provides network data that meets the exact requirements of the research and analysis, as used in the author's projects [2][16][17]. Different technologies and different protocols require different data processing. Most of the available solutions are for Ethernet based networks. Available tools that enable real-time analysis of PON networks, *International Telecommunication Union – Telecommunication* (ITU-T) G.984 [18] or ITU-T G.987 [19] are limited. Such a network card used within the project can be connected to the laboratory PON network. This card enables forwarding of traffic (PON frames) from the splitter to the development server via *Peripheral Component Interconnect Express* (PCIe). The development server then contains a software frame parser and applications (modules) for traffic analysis. To obtain outputs from this system in terms of detected security problems in the optical network, reported incidents in Apache Kafka can be stored using specific message format and then the events can be processed using custom defined templates for optical networks.

## 3.4 DATASETS

**The authors' contribution:** A specific map dataset created in JSON format for geolocating *Internet Protocol* (IP) addresses [20] [21]. Definition of GPON and XGPON frames in JSON format [2] as a source for frame and traffic analysis.

---

Datasets are used among others in the development of new methods and algorithms. A dataset, in terms of network traffic dataset, is a set of measurements that has been obtained

using data acquisition tools. They can therefore contain the data itself, or be in flow or binary form. Nowadays, many datasets that have already been created can be searched, for example using Google Dataset Search. The simplest dataset format is the *Comma-Separated Values* (CSV) for tabular data. The CSV is basically used for “flat” data. For data creating a “tree”, which has multiple layers is the JSON file format used as the most common file format. Another representation of data is the lightweight SQLite format to deliver database files, which is the most used relational database today.

## 4 DATA PREPARATION AND PROCESS STREAMLINING TECHNIQUES

This chapter focuses on data preparation related to increasing the efficiency of data processing in the traffic analysis techniques. It is composed of the author’s own results and supplemented with theoretical knowledge on the subject from other sources. The first Section 4.2 providing a brief theory to data mining and clustering topic. Section 4.2.2 presents the authors’ approach for streamlining analysis processing by using genetic algorithms.

### 4.1 DATA PREPARATION

In the following section, clustering algorithms will be discussed. Therefore, it is necessary to mention, at least, some few words to the topic of data preparation. In basic, clustering determines how similar two examples are in the way of merging all of the feature data from two instances into numerical value, so the size of data must match for the feature data when they are combined. The normalization, transformation, and creation of quantiles techniques are used besides of sampling, scaling or randomization.

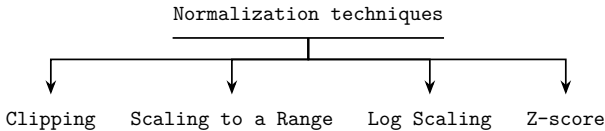


Fig. 4.1: The normalization techniques.

The normalization is used, when the data are numeric. The basic normalization techniques are shown in Figure 4.1. When the data are categorical (when features have a specific set of possible values), hashing techniques may be used to process the different length of features. In machine learning known as feature hashing techniques or hashing trick. This hashing technique is used very frequently in network analysis or network pattern analysis. The idea is very simple, to convert data into a vector of features.



## 4.2 DATA CLASSIFICATION AND CLUSTERING TECHNIQUES

**The authors' contribution:** The aforementioned subject of clustering techniques was used in the author's several projects for the preparation and classification of network traffic. Parallelization of genetic algorithms was proposed [23] [24] and used to increase the data processing efficiency of machine pre-processing. The following text and results are provided from the author's research papers.

---

Clustering is one type of unsupervised learning method used in machine learning. The objective of clustering is to group comparable groupings of the incoming data. Clustering techniques are used to classify the data obtained from multidimensional observations. The data are sorted so that the difference in the data values of the group members is close to zero. Cluster analysis is concerned with the formation of just such homogeneous units. The number of data dimensions is reduced, and one variable expresses the membership of a data unit in the cluster. The basic clustering problem can be described in “*general terms*” by providing a data matrix  $\mathbf{X}_{(m,n)}$ , where  $m$  is the number of objects and  $n$  is the number of variables. The number of clusters is denoted by  $k$ . It is a decomposition of the set of  $m$  objects in dependence on the values of  $n$  into  $k$  clusters. Only decompositions with disjunctive clusters are considered. An object must belong to only one cluster  $C_k$ . The distance for all objects is calculated. This calculation yields a square-symmetric matrix, called the association matrix. A basic overview of the clustering methods based on the approaches is shown in Figure 4.2.

Non-hierarchical, centroid (partition) based clustering includes the K-Means method, the X-Means method, and the K-Medoids method. Next, the principle of the K-Means method is presented. This method is among the methods mentioned above as their basis.

The basic K-Means algorithms randomly partitions data into  $k$  clusters. It is determined by the  $k$  centroids<sup>7</sup>  $c_k$  using the concept of the average distance in the cluster. Each cluster object and its distance to the centroid are evaluated using a distance metric. If it is closer to another one, it is relocated, and the centroids are recalculated so that a new average is computed over all the elements in the cluster. The step is repeated until none of the elements cannot be relocated anymore. Mathematically, the relation  $k$  of the  $C_k$  clusters and the  $k$

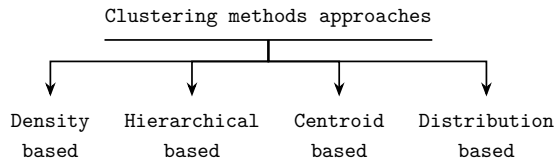


Fig. 4.2: Clustering algorithms according to their approach [22].

---

<sup>7</sup>A centroid is the center of the cluster. It is a vector containing the averages of the variables observed in the cluster.

centroids can be expressed by minimizing  $C_k$  and  $c_k$  according to the following relation 4.1 can be expressed.

$$\sum_{k=1}^k \sum_{x_n \in C_k} \|x_n - c_k\|^2 \quad (4.1)$$

The minimization problem is a hard problem to solve. The best-known solution uses Lloyd’s algorithm. Once the centroids are known, the elements are assigned according to the concept of distance according to the following relation. There are several modifications to this algorithm. Disadvantages include the fixed definition of  $k$ -clusters and the use of the Euclidean distance calculation, which is prone to distant objects. The use of Euclidean distance has drawbacks. Particularly when using high-dimensional data, the phenomenon known as the “curse of dimensionality” is affecting the result. Data get scarce as dimensionality rises, which is a concern. The parameters issue does not just apply to the issue that is discussed; it also arises with every data mining work, and the parameters have a distinct impact on every method that is utilized. Again, methods exist to validate the number of  $k$ -clusters, such as the silhouette validation method or the *Davies-Bouldin validation index (DB)*. Other clustering methods that use different distance metrics can also be used in the clustering. It depends on the format of the data, as stated before.

In network analysis, it is often necessary to compare the similarity of data strings from pre-processing. The data string can be, for example, a message hash generated or a fingerprint [25]. For this, the theory of metric space and metrics is applied. A mapping of  $\mathcal{M}^2 \rightarrow \mathbb{R}$ , where  $\mathcal{M}$  is any non-empty set, is called a metric  $\rho$ . For each  $x, y, z \in \mathcal{M}$ , the metric must meet the basic three axioms, as is identity  $\rho(x, y) = 0 \iff x = y$ , symmetry  $\rho(x, y) = \rho(y, x)$  and triangle inequality  $\rho(x, z) \geq \rho(x, y) + \rho(y, z)$ .

A standardized real and discrete metrics for different sets  $\mathcal{M}$  may be used. From the real metrics, it can be Euclidean and Manhattan. From the discrete metric, Levenshtein (edit distance is another name for this metric) and Damerau-Levenshtein metric used in combination with K-Median or *Ordering Points to Identify the Clustering Structure (OPTICS)* come into consideration, or similarity metric suitable in combination with CD-HIT clustering method. A good results were obtained by authors [25][26] with the use of the Levenshtein metric.

When using the Levenshtein metric, the least number of operations required to convert one string to another serves as “deciding factor”. The following character-related actions are acceptable: removing, adding, and swapping out characters. Levenshtein metric has been updated by the Damerau-Levenshtein metric. The same actions are permitted, plus the option to swap two consecutive characters is added. This statistic may be normalized in relation to the string’s length. The outcome is scaled to the range  $[0, 1]$ . The distance between an  $i$ -symbol prefix of string  $a$  and a  $j$ -symbol prefix of string  $b$  is known as the Damerau-Levenshtein distance between the two strings  $a$  and  $b$ .

The use case for the presented algorithms can be finding of non-usual traffic patterns. OPTICS which is a generalization of *Density-based Spatial Clustering of Applications with Noise (DBSCAN)* in combination with the Damerau-Levenshtein metric can be used to cluster

similar and non-similar network traffic patterns before being converted to fingerprints to find any anomalies. The samples are ordered in such a way that the two nearest samples always follow each other in sequence. By dividing the sequence into locations where the relative distance exceeds a predetermined threshold, the resultant clusters are identified. Separated outliers can be identified as unusual traffic detected.

To prove and evaluate the results, one of the basic methods is used to evaluate and compare univariate clustering algorithms, that is, to compare the change in cluster size for different clustering algorithms. That is, what clusters are produced by each clustering method and how they differ from each other. For this purpose, the normalized Shannon entropy  $H$  [27] can be used. The normalized Shannon entropy is defined as follows in Equation 4.2, given the proportions  $p_1, \dots, p_N$  of cells assigned to each of the  $N$  clusters.

$$\frac{H}{H_{max}} = - \sum_1^N p_i \frac{\log_2 p_i}{\log_2 N} \quad (4.2)$$

Since the actual degree of equality of cluster sizes may vary between data sets, it is useful to subtract the normalized entropy calculated from the actual distribution to obtain a final performance index. To evaluate how well inferred clusters recover the true of sub-populations and to evaluate the stability of clusters, *Hubert-Arabie Adjusted Rand Index* (ARI) [28] is one of the useful method. To evaluate, whether points are clustered well and separated well, *Silhouette Coefficient* (SC) metric can be used [29]. To obtain the silhouette score, the separation of the clusters based on the distances between and within the clusters is measured. For each sample, the mean intra-cluster distance  $A$  and the mean distance of the nearest cluster  $B$  is calculated. Then the silhouette coefficient for the sample is  $(B - A) / \max(A, B)$  with a result between 0 and 1 looking for a higher value.

When working mainly with clustering, visual examination of clusters is one of the good practices to evaluate them, mainly when they are unsupervised. To better visualize clusters, there are two very popular methods used for dimensionality reduction such as *Principal Component Analysis* (PCA) or *t-distributed Stochastic Neighbor Embedding* (t-SNE).

□ *A part of the author's research was presented here. A combination of clustering method and metric used to analyze and cluster traffic patterns. There is no simple answer to which clustering algorithm to use. It always depends on the type of data. Data sets can have many entries, but not all clustering algorithms scale well. For example, if computing the similarity between all pairs, the runtime increases as the square of the number of an entry  $n$ , thus it has a complexity of  $\mathcal{O}(n^2)$ , so such an algorithm is not practical when there are entries in millions.*

#### 4.2.1 The Role of Evolutionary Algorithms

The motivation to use the principles of evolutionary algorithms in the problem of data clustering is mainly their application to practical problems that cannot be solved with other methods. Evolutionary algorithms represent suitable techniques for solving complex optimization

problems and can achieve better results than linear methods. The most commonly used evolutionary algorithms are genetic algorithms and combined evolutionary strategies. Evolutionary algorithms are used for both single-criteria optimization and multi-criteria optimization, also called multi-objective or multi-objective optimization [30].

#### 4.2.2 Process Streamlining by Using Genetic Algorithms

Modern genetic algorithms are evolutionary algorithms that are derived from natural laws and phenomena. By their very nature, genetic algorithms lend themselves to parallel processing, which boosts performance and promotes optimization. Algorithm parallelization is a valuable technique for making an algorithm be more efficient and faster. GAs are very suitable for a large set of problems, but some of them require a more significant amount of time, and therefore, GAs became unusable for them. The most time-consuming operation within GAs is the fitness function<sup>8</sup> evaluation. This function is performed for each individual (solution) and is independent of the others. This makes it suitable for parallel processing. Genetic operators: mutation and crossover can work in isolation as they act on one or two individuals. These operators are usually much more straightforward than fitness functions but can consume more time than calculating a fitness function depending on the crossover/mutation operation. Communication is also a problem for another genetic operator, selection, which often needs information about the entire population. Therefore, the following is concerned only with parallelizing the fitness function [24]. For parallel processing, there are several basic models. They are the Hierarchical Model, the Multi-Population Coarse-Grained Model, the Global One-Population Primary-Secondary (Master-Slave<sup>9</sup>) Model, and the architecture One-Population Fine-Grained Model. The Primary-Secondary model works with only one population. Further processing, if the data of the last not yet evaluated individual in the population has already been sent, divides the Primary-Secondary GA into two types:

- **Synchronous** – in which the primary node will start sending data from the new population only after the previous one has been processed on all secondary nodes.
- **Asynchronous** – when the data of all individuals in the population have already been sent to the secondary nodes, the primary node starts sending the data from the new population to the free nodes.

Different node topologies and their neighborhoods imply different GA behaviors [23]. Frequently used topologies are one and two-dimensional grids, which are also often used to deploy computational elements of parallel computers. The Coarse-Grained model is described in [31]. The factors SpeedUp, Efficiency, and Scaling are frequently used to evaluate the advantages of parallelization [31]. The SpeedUp parameter  $S$  is given by Equation 4.3.

$$S = \frac{T_S}{T_P}, \quad (4.3)$$

---

<sup>8</sup>In general, evolutionary techniques use an objective function  $f$ , also called fitness function, to evaluate the best solution (individual)

<sup>9</sup>Master-Slave indication is replaced by Primary-Secondary.

where  $T_S$  is the computation time for the serial algorithm and the  $T_p$  is the computation time for the parallel algorithm. The efficiency parameter  $E$  is given by Equation 4.4.

$$E = \frac{S}{p}, \quad (4.4)$$

where  $p$  is the number of processing units and corresponds to the number of processes. The scaling parameter detects the loss of algorithm performance as the number of processing units and the difficulty of the calculation increase.

The contribution of parallelization can generally also be evaluated by comparing the time taken by each compared model to find the best solution. The efficiency of the algorithm itself can be evaluated by comparing the number of iterations. As an example, in the author's research [24], test results of the measurements of the time required to find the correct solution have been provided. As with the progression of the number of iterations, the similarity between the progressions of the Serial and Primary-Secondary models can be noticed, as well as the similarity between the Fine-Grained and Coarse-Grained models. It is assumed that the Fine-Grained model is faster for a sufficient number of processors and independent of population size.

□ *A part of the author's research was presented here. Parallelization applied to data clustering processes using genetic algorithms is practically applicable in network data processing and leads to an increase in the computational speed of algorithms.*

## 5 SELECTED TRAFFIC ANALYSIS TECHNIQUES

The previous chapters covered data processing. This chapter presents selected techniques that combine the various approaches and techniques in practice. It is divided into three main topics. Section 5.1 discuss the network localization techniques and provides an overview of the technique proposed by the author, Section 5.2 provides an overview of the author's contribution to this topic which is followed by Caption 6, related to the network data analysis and reporting.

**The contribution:** The sub-objective of the research was directed toward early detection of attacks and analysis of traffic behavior. In particular, the research started by focusing first on the possibilities of improving the localization analysis of stations or computer nodes in the network based on the knowledge of the general boundary of their occurrence and on the possibilities of botnet traffic detection and analysis. Second, the research continued by inventing a method and algorithms' combination for traffic similarity observation using a genetic algorithm and clustering. The articles related to this research were published in [20, 21, 32, 33, 34] and the following statements are from the results of the research carried out. The sub-part of the research is focused on the potential use of artificial intelligence as anomaly detection tools and integration into tools used for event logging. Ongoing current research focuses on the use and implementation of algorithms to enhance security in smart grids and the specific use of federated learning.

## 5.1 LOCATION ACCURACY ANALYSIS

Obtaining findings about the behavior and location of network nodes in real time is crucial for modern network security solutions. For example for *Web Application Firewall* (WAF) services and *IP Intelligence* (IPI) when using DPI.

The well-known regional Internet registries, which provide IP addresses to businesses in their respective service regions and are the main sources of data on IP addresses, are the African Network Information Centre (AfriNIC), American Registry for Internet Numbers (ARIN), Asia-Pacific Network Information Centre (APNIC), Latin American and Caribbean Internet Address Registry (LACNIC), and RIPE Network Coordination Centre (RIPE NCC)<sup>10</sup>.

There are many data mining techniques in this field, using the registry entries or specialized databases that maintain information related to IP latitude and longitude positions. The problem is with the accuracy techniques, which are not many.

Continuing research compared different algorithms and a new approach using an intersection method, the inclusion of a point in a polygon has been proposed by the author. Such an algorithm based on the crossing number method proved to be accurate also for points lying close on a cluster boundary, but with a specific issue with very close lying points. To avoid the so-called “degradation point”, starting point verification method should be used in combination with the algorithm.

Determining the inclusion of a point  $P$  in a 2D planar polygon is a geometric problem. In the research, author used the crossing number method. Within this method, the number of times a ray starting from the point  $P$  crosses the polygon boundary edges is counted. The point is outside when the number is even; otherwise, when it is odd, the point is inside [35].

The algorithm is composed of the following steps, shown in Fig. 5.1:

- In the first step, it creates a horizontal line on the right side of every point  $P$  and extends it to a defined value expressing infinity  $i$ .
- In the second step, it counts the number of times the line intersects with polygon edges.
- The conditions are determined in the following steps; a point  $P$  is inside a polygon if

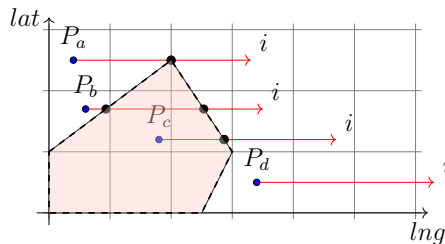


Fig. 5.1: The principle of polygon intersection.

<sup>10</sup>RIPE NCC Database: <https://www.ripe.net/manage-ips-and-asns/db/tools/geolocation-in-the-ripe-database>

<https://www.ripe.net/manage-ips-and-asns/db/tools/>

either count of intersections is odd or the point  $P$  lies on an edge of the polygon. If none of the conditions is true, then the point  $P$  lies outside.

The algorithm returns the Boolean value *true* if a point  $P$  lies on the border, or if the point  $P$  has the same value as one of the vertices of the given polygon. To do so, after the algorithm checks if the line from  $P$  to the extreme intersects, it continues to check whether the  $P$  is colinear and if the point  $P$  lies on the current side of the polygon. If it does not lie, it returns the value *false*, else *true*. The algorithm based on the crossing number method proved to also be accurate for points lying close to the map boundary. To avoid the degradation point, a verification of the starting point needs to be taken into account. This verification occurs because the boundary is also considered part of the map. The complexity of it in the worst case includes the possibility that a point  $P$  may intersect all  $N$  edges of the polygon, and  $\Omega(N)$  time is necessary in the worst case [36]. Using this algorithm and approach, it is possible to unambiguously confirm or deny the affiliation of an IP address with the map base.

## 5.2 NETWORK TRAFFIC BEHAVIOR ANALYSIS

The second aspect of network traffic behavioral analysis is to get an understanding of the specific behavior of individual network nodes, not just their specific relative or real location. There have been several methods employed up until now for identifying malicious traffic (some kind of anomaly behavior).

From this point of view, Host-based detection, Network-level-based detection, and Graph-Theory-based detection category may be used to group existing approaches. The author of this theses and his research has been generally concerned in Network-level and Graph-Theory combined-based detection. The behavioral representation of a specific network traffic connection taking into account aspects such as life time of the traffic connection, the size of the flow, and duration of the traffic flow. The data source can be a network flow protocol or, for example, a pcap file.

The related research and technique presented here focused on the possibility of comparing data flows and detecting, for example, the source of ransomware<sup>11</sup> propagation from several sources in different time sequences without using DPI, but only NetFlow protocol. The proposed technique is based on the hypothesis that each specific traffic has some unmistakable property such as periodicity in time, shape or error (traffic similarity observation) and, recently, a statistical method based on survival analysis<sup>12</sup> has been used in combination with NetFlow information. Survival analysis was originally developed to measure the lifespan of individuals. This analysis can be applied to any process duration. The survival curves, in which the traffic is transformed using the survival function based on the traffic properties can then be compared for similarity using one of the clustering techniques [34][32][20].

---

<sup>11</sup>Ransomware is a type of malware that threatens to expose a victim's personal data or permanently block access to it unless a *ransom* is paid.

<sup>12</sup>Time to event information is often the subject of survival analysis. It includes methods for positive-valued random variables in the broadest sense. Survival data are typically censored rather than fully observed.

In the case of finding some unwanted traffic, for example ransomware, it does not decide what the curve looks like. Ransomware generates specific traffic to C&C servers in time sequences and with a certain type of traffic. If traffic capture is done on multiple probes and traffic is converted to curves, the curves that are most numerous and have the most similar waveform can be filtered out using clustering or other method comparing similarity<sup>13</sup>.

The compliance tests of the survival functions are used to compare two survival curves. There are many types of these tests, each of them has optimal properties for different situations. In this case, different method as is a clustering or genetic algorithm using the Euclidean distance and the Davies–Bouldin validity index to evaluate an individual (fitness function) may provide better results[34] in term of scalability. The proposed GA algorithm to cluster the lifelines works only with centroids as individuals.

### 5.3 MACHINE LEARNING ANALYSIS

An area of computer science known as AI is concerned with the development of systems that can solve complex problems such as recognition or classification, for example, in the fields of image processing or the processing of written or spoken language, or planning or control based on the processing of large volumes of data. The term Artificial Intelligence therefore covers a set of techniques and approaches such as Machine Learning, Neural Networks, Bayesian Networks, Evolutionary Algorithms and others. Machine learning has already been used for a wide range of tasks and is particularly crucial for any application that requires the collection, analysis, and act on large data sets. The Figure 5.2 shows the possible division of machine learning algorithms based on the goal what to achieve with data. There are many more algorithms that are used depending on what the intent is. For anomaly detection analysis, i.e. anomalies

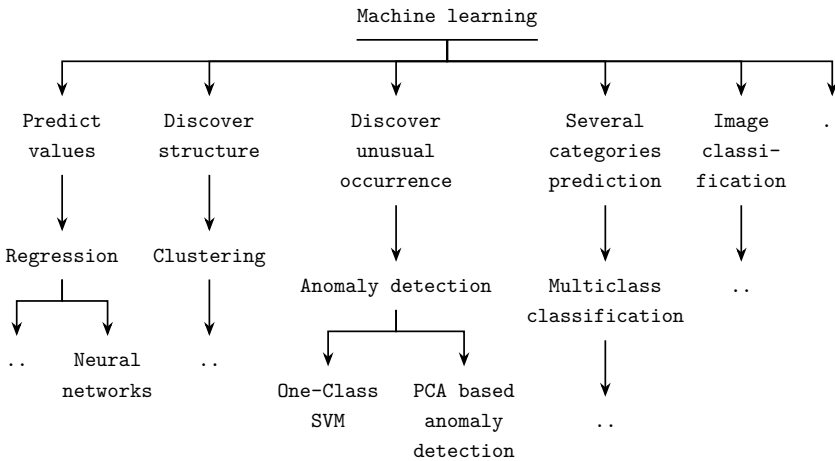


Fig. 5.2: Machine learning algorithms according to their approach [37].

<sup>13</sup><https://nsr.utko.feec.vutbr.cz/VI2VS428.php>



in traffic, algorithms like One-Class *Support Vector Machine* (SVM) or PCA Based algorithms can be used. Further, also prediction algorithms such as some type of neural network. There are three main categories of machine learning: supervised learning, unsupervised learning, and reinforcement learning [37].

The following research has focused on the method to evaluate downstream traffic behavior in the GPON ITU-T G.984 [18] transmission protocol on the sequence of *Physical Layer Operation Administration and Maintenance* (PLOAM) message ID<sup>14</sup>s during the activation process [38, 39]. The processes that an inactive ONU takes to connect or reconnect to a PON are described by the activation procedure. Three stages make up the activation process: parameter learning, serial number acquisition, and ranging. The ONU obtains the operational parameters required for the upstream transmission during the learning parameter phase. The OLT finds a new ONU (by serial number) and gives it an ONU identifier (ONU-ID) during the serial number acquisition phase. There are several states and messages used for negotiation during the activation process defined by the Recommendation ITU-T G.984.3 [40].

The objective of the solution was to use and test a machine learning-based solution for the verification of the activation phase to determine if the device under test complies with the defined recommendations and the standard, specifically, the PLOAM messages. Whether the concurrence of PLOAM messages is in accordance with the standard and whether their content is in accordance with the standard [41]. Traffic data from the network was collected using the FPGA network card, and the frames, see Figure 5.3, were parsed using a software parser [42]. PLOAMd fields were extracted from each frame. The analysis of the frames taking into account two areas of the audit:

- **Syntax verification** – examining the message to see if it complies with the standard, verification of each field content in GPON header, whether it is similar to patterns from baseline traffic or not. For the testing, One-Class SVM and AutoEncoder<sup>15</sup> have been selected. This is because there are typically few outcomes, and more significantly, these outcomes typically do not have any further structure and binary classification is in use.
- **Sequence verification** – controls the continuity of individual messages and the content of the respective fields between messages. The analysis of patterns verifies whether the analyzed protocol uses the same message in the same order and with similar content.

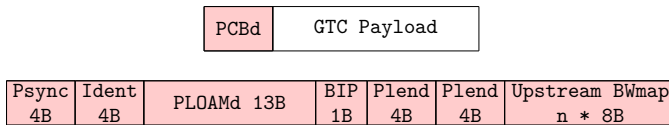


Fig. 5.3: GPON frame in downstream direction.

<sup>14</sup>Identification number, abbreviation for identification

<sup>15</sup>AutoEncoder is not a classifier, but it can be used as a layer before a classification layers. The reason to use AutoEncoder is to get a better representation of inputs, it is a dimensionality reduction technique like PCA, but it is a nonlinear dimensionality reduction.

For testing, the *Long Short Term Memory* (LSTM) network and AutoEncoder are used. This is because in the sequence verification is required to retain the information gained from previous data, and the LSTM is capable of storing processed information about the longer sequence of data.

Outlier detection using One-Class SVM classifies normal and abnormal GPON frames well, but it uses an approximate function unable to learn the importance of frame field usage. The LSTM model can distinguish time sequences. The disadvantage is that the model needs to be learned with corrupted or improper communication samples, which are not available in all cases. In combination with the AutoEncoders, the AutoEncoders prove their outlier detection capabilities and make a great alternative to the LSTM model, especially due to unsupervised learning [41].

## 5.4 ANALYSIS USING ARTIFICIAL IMMUNE SYSTEMS

AIS are adaptive systems inspired by the biological immune system. Cells, such as neutrophils, macrophages and dendritic cells are present in innate immunity. The *T* and *B* lymphocytes play a major role in adaptive immunity. In recognition, the presence of receptors that have the ability to recognize and capture a pattern is important to us. This part of the cell is called antibody. Each antibody is capable of recognizing a particular pattern. This pattern is called antigen. The principle of antigen-based pattern recognition could be compared to a key and a lock. There are also other antigens called self-antigens or self-items. These are proteins which are naturally occurring. The antibody is only capable of recognizing and binding to one particular antigen. The adaptive part of the immune system has memory. In artificial immune systems, this memory represents a trade-off between memory requirements and the rate of immune response. The use of AIS is very broad, a basic overview is shown in Figure 5.4, from classification, robotics to bioinformatics.

There are basic algorithms used in artificial immune systems [43] as are Negative selection algorithm (Forrest, Gasputa, Kim, -...), Clonal selection algorithm (Clonalg – De Castro, B-cell – Kelsley), Immune Network algorithm (models) – continuous models (Farmel, Jerne, discrete models – RAIN (Timmis), or AiNET (De Castro). The principles of the respective algorithms are well described by [43]. The techniques and algorithms used within the AIS

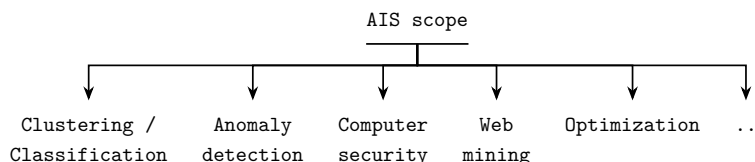


Fig. 5.4: The artificial immune system use.

are affinity<sup>16</sup> functions related to distance or similarity (Hamming, Euclidean, Manhattan), the Bone-marrow algorithms or somatic hyper-mutation. The operators used include affinity evaluation operator, operator assessment of individual levels, incentives meter calculate son, immune selection operator, clone (individual multiplication) operator, mutation operator, operator and population suppression cloning Refresh operator [44]. The **Clonal selection** algorithm has the following layers:

- Initialization – creates a random population  $P$ .
- For each antigenic pattern in a data set  $S$  do:
  - Affinity evaluation – present it to the population  $P$  and determination of affinity with each element in the  $P$ .
  - Clonal selection and expansion –  $n$  highest affinity elements of  $P$  are selected and clones are generated proportional to their affinity with the antigen, a higher affinity creates more clones.
  - Affinity maturation:
    - Each clone mutation – when high affinity, low mutation rate, and vice versa.
    - Mutated individuals are added to the population  $P$ .
    - Best individuals re-selection, kept as memory  $m$  of the antigen.
  - Metadynamics –  $n$  individuals with low affinity are replaced by randomly generated new  $n$ .
- Cycle – the second step is repeated until a certain stopping criterion is met during the cycle.

The **negative selection** algorithm, shown in Figure 5.5, has the following layers:

- Self-definition as normal pattern:
  - A set of equal size of the pattern sequence.
  - The set is presented as a multi-set  $S$  of string of length  $l$  over a finite alphabet.
- A set of  $R$  detectors is generated, each of which fails to match any string in  $S$ .
- The changes of  $S$  are observed by testing the  $R$  matching the  $S$ , and if any  $R$  matches, the non-self is detected.

When working with AIS, the optimization of the high density of individuals is constrained to ensure the diversity of individuals in a given solution. The antibody-antigen affinity density is defined by Equation 5.1 as the SMA of antibody-antigen affinity.

$$den(Ab_i) = \frac{1}{m} \sum_{i=1}^m aff(Ab_i, Ag_i), \quad (5.1)$$

where  $Ab_i$  is the antibody of the  $i$ -species, the size of the population is defined by  $m$  and  $aff$  is the antibody-antigen affinity for  $i$ .

The antibody-antigen affinity  $aff$  is mainly based on the method of calculating the Euclidean distance, Hamming distance or the information entropy. If the  $aff$  is smaller, a stronger

---

<sup>16</sup>Affinity – Strength of ligand (molecule) binding to its receptor.



Fig. 5.5: The negative selection algorithm.

affinity is encountered, thus, the stronger ability of the antibody to capture and kill an antigen, in the case of network analysis, do the detection. Problems like function approximation and optimization are being solved using a variety of machine-learning techniques.

Among these, *Artificial Neural Network* (ANN) for non-linear function approximation and the use of genetic algorithms to find the optimum (maximum or minimum) of a function. The starting collection of candidate solutions for the GA or initial weight vectors for the ANN must be defined for both procedures [43].

Antibodies can be modeled as a string of bits of length  $l$ . As a binding of antibody to an unknown element, a string match of bits on the antibody and a string of bits on the antigen (unknown element) is considered. In order for the antibody to bind to the antigen (the affinity), it is necessary that the chains equalize each other. But this is hardly a feasible requirement, so the principles of string similarity are used. These include the Hamming distance, the Euclidean distance, or the Manhattan distance, as mentioned above. It can be a binary representation, continuous as numeric, or categorical. Assume the general case  $Ab = \langle Ab_1, Ab_2, \dots Ab_i \rangle$  and  $Ag = \langle Ag_1, Ag_2, \dots Ag_i \rangle$ . The binary representation typically employs the Hamming rule, and the numeric (real or integer) representation typically Euclidean. When using the principle of AIS in clustering base algorithms, data items are seen as antigens and clusters as antibodies. The process by which the immune system repeatedly creates antibodies to detect the antigen and finally creates the best antibody that can capture the antigen is analogous to the clustering of the data items [45] and clustering K-Means, K-Nearest based algorithms are used.

#### 5.4.1 Mapping selected problem into the AIS

Applying the same problem as in Caption 5.3 for the activation process, in the case of an recognition system, antigens correspond to the GPON frames whose contents may contain some unidentified field, for example, unspecified PLOAM fields and its content by the standard. Self-antigens match the known parts of the GPON frame. An antibody presents a bit pattern that

matches a potential not known part of the specified GPON frame. A lymphocyte represents two or more detectors. The cell apoptosis<sup>17</sup> is modeled using the negative selection algorithm. The steps of the algorithm are as follows [11]:

- Input: self-set (known patterns). Output: New set  $P$  of detectors.
- Initialization of the empty set  $P$  of detectors (memory cells) and determination of the affinity boundary (maximum similarity between the detector and some element of self-set).
- Creating a random detector.
- Determining the affinity of the detector step by step for all elements of the self-set and as the detector affinity will be considered the highest one.
- If the solution has less affinity than the determined affinity limit, add the detector to the  $P$  set.
- If the  $P$  set is large enough, terminate the algorithm. Continue with the detection.
- Input: Set  $P$ , Data  $D$ .
- Count affinity of Set  $P_i$ , Data  $D_i$ .
- If affinity and matches  $D_i, P_i$  – detection.

The operators used are in this case hard-coded; thus, they are not included in the algorithms' steps. The optimization function is done by using GA, thus, this is a mix of genetic and K-Nearest algorithm. When comparing ML Syntax model and AIS algorithms, the AIS algorithm showed a higher success rate when recognizing modified strings of data, in this case PLOAM messages, but only when the ratio of modified strings was lower than the correct messages. However, the performance aspect and the processing time must be taken into account. When the set of correct messages is lower than the possible set of non-self, it fails. The negative selection algorithm serves to detect detectors that can only recognize foreign elements. Thus, this algorithm removes the elements that recognize the known element. It is used where the known set is much larger than its complement.

## 6 NETWORK DATA ANALYSIS AND REPORTING

This chapter presents the current data monitoring situation from a legislative perspective and from the perspective of active elements. Then the principle of the possibility of analysis and monitoring of passive optical networks is explained, which is one of the author's contributions in this topic.

**The contribution:** The research was directed toward the design and development of an active network element and the implementation of algorithms such that would enable efficient analysis of transmitted data structures in a real-time optical access and distribution network on the GPON ITU-T G.984 [18] transmission protocol, and the continuous research on the XGPON ITU-T G.987.1 [19] [16] [2]. This research also involved the upgrade of existing data processing algorithms and algorithms used in the first sub-objective research. During the research, it was

---

<sup>17</sup>The death of cells that occurs as a normal and controlled part of an organism's growth or development

also necessary to resolve issues related to the large amount of data processing for the analysis itself and the preparation of data for the individual processes [46]. The obvious solution was to introduce the parallelization of computational processes and algorithms [24]. Part of the research and development involved designing a system that would be able to analyze traffic in passive optical networks in real time. Such a system was lacking on the market at the time the research was initiated. Real-time traffic analysis is provided mostly at the higher Ethernet layers. However, the intent of the research conducted was to analyze real-time management traffic to detect deviations from the standard, errors, or unsolicited traffic directly at the passive optical network layer.

---

Data structure analysis is nowadays equally important. In order to further strengthen the resilience and incident response capabilities of the public and commercial sectors as well as the *European Union* (EU) as a whole, the Council and the European Parliament agreed on steps for a high common level of cybersecurity across the Union. The present *Network and Information Systems* (NIS) directive regulation on the security of network and information systems will be replaced once the new directive, known as NIS2 [47]. A range of tools are used by security teams to investigate and mitigate breaches. The selection of certain tools is left up to the discretion of the many separate teams; there is no uniform standard that dictates its use. Teams are given particular advice and guidelines by organizations like *European Union Agency for Cybersecurity* (ENISA) and *National Institute of Standards and Technology* (NIST). *Intrusion Detection Systems* (IDS), *Intrusion Prevention Systems* (IPS), *Security Information and Event Management* (SIEM), and *Security Orchestration, Automation and Response* (SOAR) are a few broad categories of the technologies used.

The following author's research and work have been motivated by the non existence of a solution for GPON and XG-PON networks related to the verification of standard and automatic reporting features. The security features specified by the standards for the PON network are based on the assumptions that eavesdropping on the signal is not a trivial task. Uplink transmission in the two-direction communication of the PON network is therefore considered secure. However, this may not always be true in a real environment. Splitters<sup>18</sup> are commonly installed, for example, in basements where they are not further secured. If some of the ports on the splitter are free, an attacker can easily connect to the network; otherwise, it is enough to disconnect a legitimate user from the network and connect an attacker's device known as Rogue ONU. Several known types of security risks should not be ignored within PON. It also presents some other of the security weaknesses of GPON. Modified ONUs represent one of the most significant security risks in the PON network. For example, they can be used for attacks such as *Theft of Service* (TOS), Masquerade, or Reply Attack. *Denial of Service* (DoS) attacks make a network service unavailable to legitimate users. The blocking of upstream communication occurs when an ONU transmits outside of its allocated time slots. The root cause of a

---

<sup>18</sup>The passive optical splitter can split, or separate, an incident light beam into several light beams at a certain ratio.

DoS attack also can be a hardware or software malfunction of an ONU. However, an attacker can deliberately modify an ONU to transmit continuously on a given wavelength and with sufficient transmit power to block the communication of other ONUs. The attack can be realized with a sufficiently powerful laser beam source [48]. Specific exploits are available in the Exploit Database [49].

As part of the research, the exploitation of exposed ONUs on the Internet, as well as their “Web App” API interfaces, and how it would be possible to paralyze an internal network [50] has been investigated. The research focused on PLOAM messages used in the management communication between OLT and ONU. These messages are used to transmit control and monitoring instructions between the OLT and the ONU, but in the theory, it would be possible the messages can carry unwanted instruction. It has been found, that some vendors do not follow the standard, and thus, undefined messages are present in these PLOAM messages.

The individual security tools work with *Indicators of Compromise* (IoC), forensic data discovered during network or system monitoring indicating potential intrusion or malicious activity. In general, these can be, for example, IP addresses, malware signatures, domain names, malware file hashes, and more [51]. Another example can be unusual activities such as data found in system log entries, unusual network traffic, bundles of data in the wrong place, etc. According to the research findings, most systems are designed for IP networks, but little attention is paid to the optical access network. The IoCs are *not defined for PON networks*. For example, based on the detection of non-standard PLOAM messages, these particular messages could be considered as IoCs. Several possible IoC to be used for PON system automation and monitoring has been proposed with this research.

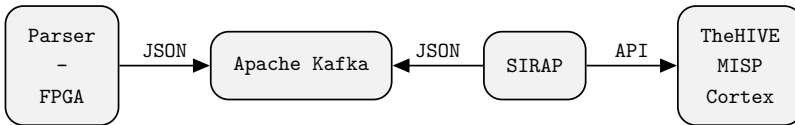


Fig. 6.1: The functional diagram of the reporting solution.

*Security Incident Response Automation for xPON* (SIRAP), shown in Figure 6.1, consists of an API and parts, published in [2], that provide a collection and analysis of reports from PON networks. The SIRAP also uses modules that connect it to existing incident reporting tools using specific IoC identifiers for PON. The SIRAP is a middleware between the rest of the modules processing data and tools used for reporting.

Specific developed FPGA card [16] for data acquisition is connected to capture data on the PON layer and transfers frames from the optical network (downlink and uplink direction) to the server’s *Direct Memory Access* (DMA). The frame parser [42] is a C# application parsing the traffic from the FPGA, creating a JSON with parsed frames and sends them for further processing to the Apache Kafka. The Apache Kafka is used to scale the solution when working with high speed data, buffering the messages. The data analysis module included in the SIRAP

is based on TensorFlow, meaning that it is a TensorFlow detector. Once the traffic is analyzed and if unspecified or unsolicited traffic is found, a report is created according to a template specific to TheHIVE<sup>19</sup> and sent via API to TheHIVE. So far, the following report types were defined:

- **PLOAMd Anomaly** – notification of anomaly detection in a PLOAMd message, i.e., a deviation from the messages specified by the standard.
- **Activation process anomaly** – notification of anomaly detection in activation process.
- **Non-standard Frame structure** – notification of anomaly detection in PON Transmission Convergence Layer frames.
- **ONU Management and Control Interface (OMCI) Anomaly** – notification of anomaly detection in OMCI messages or the *ONU Management and Control Channel (OMCC)* activation process channel.
- **Non-specified error** – the empty report, used as a template for the specification of a new report type.

This solution allows both monitoring of the passive optical network and defining custom templates for incident reporting. Reports can be created using the defined API.

## 7 VERIFICATION AND INTERPRETATION OF RESULTS

The previous chapters have partially presented methods for verifying results related to each selected traffic analysis technique, such as for clustering algorithms. This chapter provides a general summary overview dealing with the verification and interpretation of results obtained from a network analysis. The purpose of this chapter is to summarize in general terms the requirements for the last part of the traffic analysis presented in Chapter 2, namely the general procedure for verifying the results. The following pattern of approach is more a general conception of the author of the thesis than an established paradigm and it is open for a discussion.

### 7.1 GENERAL APPROACH

The general verification model can serve well as a pattern. A pattern [52] is a way of regulating a group of different accesses using a mechanism that is defined in a particular environment. The goal of the pattern is not to eliminate the weaknesses of the approach; rather, it is to minimize or mitigate the error of the approach. When working on interpreting the results of traffic analysis, it is useful to have basic access parameters set. First and foremost, it is necessary to ask whether the analysis performed is correct. If possible, it is advisable to compare and quantify the results with some known data. Accuracy and precision are two measures of observational error that can be observed. Precision in meaning means how close a given set of measurements or observations is to its true value, while accuracy means how

---

<sup>19</sup>TheHIVE <https://thehive-project.org/index.html>



close the measurements are to each other. The accuracy is defined by ISO 5725-1:1994/Cor 1:1998<sup>20</sup>. In previous chapters, some techniques for such measurement have been presented for specific purposes. If the measurements do not conform, it is advisable to review the procedure, the values used, and the semantics of the code. An example is the silhouette measurement for the research paper in Section 4.2. Another example, if a GA were required to be rated, here again the situation is perceived differently. In practice, empirical tuning is widely used so far. The best tuning of the actual GA parameters is determined by extensive experiments on their performance using simulation. Either self tuning or hierarchical tuning is used, where another GA manipulates the parameters of the measured GA and tracks the best performance.

Therefore, if the results are satisfactory, it is possible to proceed to step three. If they are still not, it is advisable to go back to step two, and be interested in the data itself. It can be helpful to visually inspect the results by displaying them using subsidiary graphs. For now, there is nothing that completely replaces the experience acquired, or at least some general idea of the result. The quality control should look at what data is available and whether it contains unknown values. For example, whether they contain *Not a Number* (NaN) or null values. The possible outcomes and values are limitless, bases on what tool and type of the measurement is used. It's also good to find out what the data layout is. Whether the data is from a normal (Gaussian) or power law or different distribution. Q-Q (Quantiles Quantiles) [53] plots, for example, can be used to determine the distribution. The data needs to be normalized and scaled based on the distribution they have. The data preparation techniques are out of the scope of this theses, some of them were presented in Chapter 4. The choice of technique depends on the type of data. When the data are normalized and scaled, analysis algorithms will calculate more precisely. Once the data analysis correction has been performed, it is advisable, if possible, to subject the results of the analysis to validation and interpret the results appropriately. In the case of evaluating the optimal solution or some other optimality, hypotheses are used. In terms of statistics, the optimality criterion provides some measure of the fit of the data to the hypothesis and helps in the selection of the model or procedure used. For example, in traffic analysis practice, it is determined how optimal an algorithm is. For example, in machine learning, overfitting and underfitting are distinguished<sup>21</sup>. These processes should reflect the current situation, and the results of the analysis should be subject to repeated inspection. Repeated verification provides a means to ensure that traffic analysis is effective and accurate.

## CONCLUSION

This habilitation thesis aimed to present a summary of the author's work and research in the field of network traffic data analysis. The field of networking technology is evolving at an unstoppable pace, and there is a need to respond to this evolution in a very flexible manner. Therefore, the different aspects of the thesis have been presented to reflect at least the current

---

<sup>20</sup><https://www.iso.org/standard/29779.html>

<sup>21</sup><https://www.ibm.com/cloud/learn/overfitting>

topics. The content of the thesis was divided into two main areas, namely theoretical and practical, which dealt with selected technologies and techniques for network traffic analysis that the author has worked with both in his projects and collaborative research activities. They also draw on experience in commercial and academic environments.

The thesis consists of seven chapters in total. The first chapter stated the objectives of the thesis and the author's own contribution and organization of the thesis. The next chapter dealt with the current knowledge in the field of traffic analysis. A generic view of the problem and of the various traffic analysis techniques was presented, with an indication of the basic mathematical relationships later in the text for each technique. Furthermore, the chapters were structured according to the generic view of traffic analysis. First, the principles of data sources for traffic analysis were explained, then how the data are processed, and then in which techniques they are used were presented. These were mainly selected techniques covering the author's contributions and publications in the field. Thus, it was not an exhaustive list of techniques in use today, but a view of some of them was offered. In particular, this included the author's proposed technique for making data localization more efficient, as well as the author's proposed technique and algorithm for analyzing traffic behavior based on the clustering algorithm and the use of the survival algorithm. Another technique discussed the use and aspects of artificial intelligence, of which the less-used traffic analysis technique using the artificial immune system is represented here. Subsequently, a practical example and the issue of monitoring and its possible solution for passive optical networks were presented. The whole work is designed to have a contribution to education as well. Thus, each chapter contained a part of the theory on the problem and a part dedicated to the author's work.

In general, there is no one-size-fits-all instruction on how to work with data in traffic analysis. The author's recommendation is to work with intuition and also use approaches used in other research fields, such as in the medical or nature field, and to use own imagination to the fullest. This aspect was reflected in the results of the author's work.

The thesis presents the results of the author's work achieved from the completion of his Ph.D. studies in 2017 to the present. Further, ongoing research in this area is mainly focused on improving existing methods of traffic behavior detection and their other applications, especially within the project "Data backup and storage system with integrated active protection against cyber threats", and "Android federated learning framework for emergency management applications". The first project focuses on creating a solution for the early detection of ransomware before backing up virtual machines. The second project focuses on developing a framework based on federated learning for secure data transfer and decentralized learning for crisis management applications.

The author of this thesis is the author or co-author of 30 articles and 29 conference papers, most of them indexed in Web Of Science, Scopus or with an impact factor, with a total of over 160 citations at the time of writing this thesis.

## BIBLIOGRAPHY

- [1] *Network Research Group*. [Online; accessed 2022-09-09]. URL: <https://nsr.utko.feec.vutbr.cz/>.
- [2] Vaclav Oujezsky, Tomas Horvath, and Martin Holik. “Security Incident Response Automation for xPON Networks”. In: *Journal of Communications Software and Systems* 18.2 (2022), pp. 144–152.
- [3] Petr Velan et al. “A survey of methods for encrypted traffic classification and analysis”. In: *International Journal of Network Management* (2015). URL: <https://doi.org/10.1002/nem.1901>.
- [4] H. Song et al. *Network Telemetry Framework*. RFC 9232. RFC Editor, May 2022.
- [5] M. Bjorklund. *The YANG 1.1 Data Modeling Language*. RFC 7950. RFC Editor, Aug. 2016.
- [6] P4.org Applications Working Group. *In-band Network Telemetry (INT) Dataplane Specification*. URL: [https://github.com/p4lang/p4-applications/blob/master/docs/INT\\_v2\\_1.pdf](https://github.com/p4lang/p4-applications/blob/master/docs/INT_v2_1.pdf) (visited on 11/28/2022).
- [7] F. Brockners, S. Bhandari, and T. Mizrahi. *Data Fields for In Situ Operations, Administration, and Maintenance (IOAM)*. RFC 9197. RFC Editor, May 2022.
- [8] Divya Somvanshi and R.D.S. Yadava. “Boosting Principal Component Analysis by Genetic Algorithm”. In: *Defence Science Journal* 4.60 (2010), p. 7. DOI: 10.1.1.902.7675. URL: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.902.7675%5C&rep=rep1%5C&type=pdf> (visited on 04/12/2017).
- [9] Wilson Rivera-Gallego. “A GENETIC ALGORITHM FOR SOLVING THE EUCLIDEAN DISTANCE MATRICES COMPLETION PROBLEM”. In: *SAC* (1998), p. 5. URL: [http://slapper.apam.columbia.edu/bib/papers/river\\_b\\_99.pdf](http://slapper.apam.columbia.edu/bib/papers/river_b_99.pdf) (visited on 04/12/2017).
- [10] Dipankar Dasgupta. “Advances in artificial immune systems”. In: *IEEE Computational Intelligence Magazine* 1.4 (2006), pp. 40–49. ISSN: 1556-603X. DOI: 10.1109/MCI.2006.329705. URL: <http://ieeexplore.ieee.org/document/4129847/> (visited on 06/24/2019).
- [11] Vaclav Oujezsky, Vladislav Skorpil, and Tomas Horvath. “Gpon frame analysis with artificial immune system”. In: *2019 11th International Congress on Ultra Modern Telecommunications and Control Systems and Workshops (ICUMT)*. IEEE, 2019, pp. 1–4.
- [12] Huawei. *6G: The next horizon white paper*. URL: <https://www.huawei.com/en/technology-insights/future-technologies/6g-white-paper> (visited on 11/27/2022).
- [13] B. Claise, B. Trammell, and P. Aitken. *Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of Flow Information*. STD 77. RFC Editor, Sept. 2013. URL: <http://www.rfc-editor.org/rfc/rfc7011.txt>.

- [14] B. Claise. *Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of IP Traffic Flow Information*. RFC 5101. RFC Editor, Jan. 2008. URL: <http://www.rfc-editor.org/rfc/rfc5101.txt>.
- [15] Vladislav Skorpil, Vaclav Oujezsky, and Ludek Palenik. “Internet of things security overview and practical demonstration”. In: *2018 10th International Congress on Ultra Modern Telecommunications and Control Systems and Workshops (ICUMT)*. IEEE. 2018, pp. 1–7.
- [16] Vaclav Oujezsky et al. “Fpga network card and system for gpon frames analysis at optical layer”. In: *2019 42nd International Conference on Telecommunications and Signal Processing (TSP)*. IEEE. 2019, pp. 19–23.
- [17] DFC Design. *DFC we MAKE Electronics...* URL: <https://www.dfcdesign.cz/cz>.
- [18] *Gigabit-capable passive optical networks (G-PON): Transmission convergence layer specification*. 2014. URL: <https://www.itu.int/rec/T-REC-G.984.3> (visited on 10/02/2021).
- [19] *G.987.1 : 10-Gigabit-capable passive optical networks (XG-PON): General requirements*. 1st ed. Geneva, Switzerland: ITU-T, 2016.
- [20] V Oujezsky and T Horvath. “Aequor Tracer–Network Analysis Application”. In: *2019 27th Telecommunications Forum (TELFOR)*. IEEE. 2019, pp. 1–4.
- [21] Vaclav Oujezsky, Tomas Horvath, and Petr Munster. “Application for Determining whether IP Addresses belong to a Map by Coordinates”. In: *2019 42nd International Conference on Telecommunications and Signal Processing (TSP)*. IEEE. 2019, pp. 14–18.
- [22] Dongkuan Xu and Yingjie Tian. “A Comprehensive Survey of Clustering Algorithms”. In: *Annals of Data Science* 2.2 (June 2015), pp. 165–193. ISSN: 2198-5812. DOI: 10.1007/s40745-015-0040-1. URL: <https://doi.org/10.1007/s40745-015-0040-1>.
- [23] V Skorpil et al. “Parallel processing of genetic algorithms in Python language”. In: *2019 Photonics & Electromagnetics Research Symposium-Spring (PIERS-Spring)*. IEEE. 2019, pp. 3727–3731.
- [24] Vladislav Skorpil and Vaclav Oujezsky. “Parallel Genetic Algorithms’ Implementation Using a Scalable Concurrent Operation in Python”. In: *Sensors* 22.6 (2022), p. 2389.
- [25] Pavel Novak and Vaclav Oujezsky. “Detection of Malicious Network Traffic Behavior Using JA3 Fingerprints”. In: *Proceedings II of the 28th Conference STUDENT EEICT 2022*. Ed. by Assoc. Prof. Vítězslav Novák. Brno: Brno University of Technology, Faculty of Electrical Engineering and Communication, 2022, pp. 194–197. ISBN: 978-80-214-6030-0. URL: [https://www.eeict.cz/eeict\\_download/archiv/sborniky/EEICT\\_2022\\_sbornik\\_2\\_v2.pdf](https://www.eeict.cz/eeict_download/archiv/sborniky/EEICT_2022_sbornik_2_v2.pdf).

- [26] Pavel Novák. “Detection of malicious network traffic behavior”. Master’s thesis. Brno: Masaryk University, Faculty of Informatics, 2022. URL: <https://is.muni.cz/th/dq2t1/>.
- [27] C. E. Shannon. “A mathematical theory of communication”. In: *The Bell System Technical Journal* 27.3 (1948), pp. 379–423. DOI: 10.1002/j.1538-7305.1948.tb01338.x.
- [28] Lawrence Hubert and Phipps Arabie. “Comparing partitions”. In: *Journal of Classification* 2.1 (Dec. 1985), pp. 193–218. ISSN: 1432-1343. DOI: 10.1007/BF01908075. URL: <https://doi.org/10.1007/BF01908075>.
- [29] R.O. Sinnott, H. Duan, and Y. Sun. “Chapter 15 - A Case Study in Big Data Analytics: Exploring Twitter Sentiment Analysis and the Weather”. In: *Big Data*. Ed. by Rajkumar Buyya, Rodrigo N. Calheiros, and Amir Vahid Dastjerdi. Morgan Kaufmann, 2016, pp. 357–388. ISBN: 978-0-12-805394-2. DOI: <https://doi.org/10.1016/B978-0-12-805394-2.00015-5>. URL: <https://www.sciencedirect.com/science/article/pii/B9780128053942000155>.
- [30] Václav Oujezský. “Converged Networks and Traffic Tomography by Using Evolutionary Algorithms”. Dissertation Thesis. Brno: Brno University of Technology. Faculty of Electrical Engineering and Communication. Department of Telecommunications, 2017. URL: <http://hdl.handle.net/11012/68296> (visited on 11/27/2022).
- [31] Vladislav Skorpil, Vaclav Oujezsky, and Martin Tuleja. “Testing of Python models of parallelized genetic algorithms”. In: *2020 43rd International Conference on Telecommunications and Signal Processing (TSP)*. IEEE. 2020, pp. 235–238.
- [32] Vaclav Oujezsky, Tomas Horvath, and Vladislav Skorpil. “Botnet C&C traffic and flow lifespans using survival analysis”. In: *International Journal of Advances in Telecommunications, Electrotechnics, Signals and Systems* 6.1 (2017), pp. 38–44.
- [33] Vaclav Oujezsky and Tomas Horvath. “Traffic analysis using netflow and python”. In: *Informatyka, Automatyka, Pomiar w Gospodarce i Ochronie Środowiska* 7.2 (2017), pp. 5–7.
- [34] Vaclav Oujezsky and Tomas Horvath. “Traffic similarity observation using a genetic algorithm and clustering”. In: *Technologies* 6.4 (2018), p. 103.
- [35] Rod Pierce. *Area of Irregular Polygons*. URL: <https://www.mathsisfun.com/geometry/area-irregular-polygons.html> (visited on 11/27/2022).
- [36] Frank L. D evai. “On the Complexity of Some Geometric Intersection Problems”. In: 1995.
- [37] Microsoft. *An introduction to the mathematics and logic behind machine learning*. URL: <https://azure.microsoft.com/cs-cz/resources/cloud-computing-dictionary/what-are-machine-learning-algorithms> (visited on 11/27/2022).

- [38] Tomas Horvath et al. “Activation Process of ONU in EPON/GPON Networks”. In: *2018 41st International Conference on Telecommunications and Signal Processing (TSP)*. IEEE. 2018, pp. 1–5.
- [39] Tomas Horvath et al. “Activation Process of ONU in EPON/GPON/XG-PON/NG-PON2 Networks”. In: *Applied Sciences* 8.10 (2018). ISSN: 2076-3417. DOI: 10.3390/app8101934. URL: <https://www.mdpi.com/2076-3417/8/10/1934>.
- [40] International Telecommunication Union (ITU). *Recommendation ITU-T G.984.3 Gigabit-Capable Passive Optical Networks (G-PON): Transmission Convergence Layer Specification*. Tech. rep. Geneva, 2014. URL: <https://www.itu.int/rec/T-REC-G.984.3-201401-I/en>.
- [41] Vaclav Oujezsky et al. “Gpon traffic analysis with tensorflow”. In: *2020 43rd International Conference on Telecommunications and Signal Processing (TSP)*. IEEE. 2020, pp. 69–72.
- [42] Michal Jurcik et al. “GPON parser for database analysis”. In: *2019 42nd International Conference on Telecommunications and Signal Processing (TSP)*. IEEE. 2019, pp. 347–350.
- [43] Leandro Nunes de Castro and Von Zuben. “Artificial Immune Systems: Part I-Basic Theory and Applications”. In: 1999.
- [44] Jing Zhang. “Artificial immune algorithm to function optimization problems”. In: *2011 IEEE 3rd International Conference on Communication Software and Networks*. 2011, pp. 667–670. DOI: 10.1109/ICCSN.2011.6014177.
- [45] Tao Liu et al. “A New Clustering Algorithm Based on Artificial Immune System”. In: *2008 Fifth International Conference on Fuzzy Systems and Knowledge Discovery*. Vol. 2. 2008, pp. 347–351. DOI: 10.1109/FSKD.2008.67.
- [46] Martin Holik et al. “Storage for Traffic from xPON Networks”. In: *2020 43rd International Conference on Telecommunications and Signal Processing (TSP)*. IEEE. 2020, pp. 77–80.
- [47] *The NIS 2 Directive*. 2022. URL: <https://www.nis-2-directive.com/> (visited on 11/19/2022).
- [48] David Gutierrez, Jinwoo Cho, and Leonid G. Kazovsky. “TDM-PON Security Issues: Upstream Encryption is Needed”. In: *OFC/NFOEC 2007 - 2007 Conference on Optical Fiber Communication and the National Fiber Optic Engineers Conference*. 2007, pp. 1–3. DOI: 10.1109/OFC.2007.4348474.
- [49] *Exploit Database - Exploits for Penetration Testers, Researchers, and Ethical Hackers*. OffSec Services Limited. 2022. URL: <https://www.exploit-db.com/> (visited on 03/26/2022).

- [50] Vaclav Oujezsky et al. “Security testing of active optical network devices”. In: *2019 42nd International Conference on Telecommunications and Signal Processing (TSP)*. IEEE. 2019, pp. 9–13.
- [51] Onur Catakoglu, Marco Balduzzi, and Davide Balzarotti. “Automatic Extraction of Indicators of Compromise for Web Applications”. In: *Proceedings of the 25th International Conference on World Wide Web*. WWW '16. Montréal, Québec, Canada: International World Wide Web Conferences Steering Committee, 2016, pp. 333–343. ISBN: 9781450341431. DOI: 10.1145/2872427.2883056. URL: <https://doi.org/10.1145/2872427.2883056>.
- [52] Eduardo B Fernandez. *Security patterns in practice. designing secure architectures using software patterns*. United Kingdom: John Wiley & Sons, Ltd., 2013. ISBN: 978-1-119-99894-5.
- [53] Paras Varshney. *Q-Q Plots Explained*. Towards Data Science, Medium. 2022. URL: <https://towardsdatascience.com/q-q-plots-explained-5aa8495426c0> (visited on 11/20/2022).

## **ABSTRAKT**

Tato habilitační práce prezentuje přehled poznatků z výzkumu zaměřeného na analýzu provozu v datových sítích. Analýza síťového provozu obecně, je nedílnou součástí soudobých síťových bezpečnostních systémů. V této práci jsou uvedeny vybrané techniky používané v analýze provozu, vlastní výzkumné závěry v této oblasti a dosažené cíle, které reflektují publikované výsledky. Práce je rozdělena do sedmi hlavních kapitol. Úvodní část stanovuje cíle práce a shrnuje současný stav vědeckého poznání problematiky a dostupných technologií. Následně je představen způsob získávání dat síťového provozu a jejich formát. Dále je v předložené práci řešena problematika zpracování dat pro další analýzu. Následuje kapitola zabývající se vybranými technikami analýzy síťového provozu a možnostmi jeho monitoringu. Závěrečná kapitola shrnuje celou práci. Všechny uváděné výsledky byly testovány nebo ověřeny v laboratorním prostředí. Vlastní text je psán tak, aby byl jak vědecky, tak pedagogicky přínosný. Práce je založena především na původním výzkumu a vývoji autora v letech po obhajobě jeho doktorské práce. Všechna prezentovaná řešení byla publikována v časopisech s impakt faktorem, indexovaných časopisech nebo prezentována na mezinárodních konferencích.

## **ABSTRACT**

This habilitation thesis presents an overview of research findings focused on traffic analysis in data networks. Network traffic analysis, in general, is an integral part of modern network security systems. This thesis presents selected techniques used in traffic data analysis, the authors' research findings in this area, and the goals achieved, which reflect the published results. The thesis is divided into seven main chapters. The opening section states the objectives of the thesis and summarises the current state of the art in the field and the available technologies. Subsequently, methods of network traffic data acquisition and data format are presented. Next, the given work addresses the question of data processing for further analysis. This is followed by a chapter dealing with selected network traffic analysis techniques and monitoring options. The final chapter summarizes the whole work. All the presented results have been tested or verified in a laboratory environment. The actual text is written to be both scientifically and pedagogically valuable. The thesis is mainly based on the author's original research and development in the years following the defense of his Ph.D. thesis. All presented solutions have been published in impact factor journals, indexed journals, or presented at international conferences.