# Infocommunications Journal

Technically Co-Sponsored by

IEEE ComSoc™
**IEEE Communications Society**

hte

IEEE
**HUNGARY SECTION**

## Indexing information

Infocommunications Journal is covered by Inspec, Compendex and Scopus.
Infocommunications Journal is also included in the Thomson Reuters – Web of ScienceTM Core Collection,
Emerging Sources Citation Index (ESCI)

# From Picosatellites to Quantum Genetic Algorithms – the latest proceedings of Infocommunications

Pal Varga

THE 2023 autumn issue of Infocommunications Journal gives a nice mixture of the latest "news" in our domain.

The article by Hussein Taha, Péter Vári and Szilvia Nagy discusses the future of the 470-694 MHz band in Europe, as the World Radiocommunication Conference 2023 (WRC-23) is set to decide on its use, balancing broadcasting services and mobile broadband needs. It analyzes various options, recommending secondary allocation to mobile service in a downlink-only mode to accommodate both broadcasting and mobile services without causing interference.

Tibor Herman and Levente Dudás introduces a new technique for identifying picosatellites and estimating their Doppler shift using passive radar methods and crosscorrelation with known transmissions in their paper. Their method enables the use of omnidirectional antennas instead of high-gain directional ones. The algorithm compensates for antenna gain differences and has been practically tested on the MRC-100 PocketQube mission. The method particularly addresses the challenge of Doppler estimation for small satellites without precise orbital data post-launch, enhancing tracking accuracy when multiple satellites are closely grouped in orbit.

The paper by Mohamed Al Amrani et al. presents a novel approach to optimizing caching and availability in UAVassisted cellular networks for 5G, employing game theory to analyze the competition among UAVs for caching and sharing revenue, leading to cost-effective strategies for energy and QoS. The study underscores the importance of UAVs as flying relays in smart city development and 5G networks, focusing on optimizing network performance through competitive game theory to achieve a Nash equilibrium that maximizes coverage and ensures fair pricing.

In their paper, Ameen Al-Azzawi and Gábor Lencse delve into the Lightweight 4over6 (lw4o6) transition technology, which facilitates the move from IPv4 to IPv6, comparing it to DS-Lite in terms of topology, functionality, and security concerns. Through a practical test-bed built with open-source software like Snabb, their study explores lw4o6's efficiency and security by simulating attack scenarios and suggesting countermeasures.

The paper by János Csatár, Péter György and Tamás Holczer investigates security weaknesses of the IEC 60870-5-104 protocol, widely used in European power systems, which lacks essential security features such as encryption and authentication. The authors developed and tested novel attack vectors in both simulated and actual environments, and in the paper highlighted potential entry points for threat actors and demonstrating how attacks can be precisely targeted to compromise telecontrol systems.

György Wersényi evaluated bone conduction and active noise-cancellation headsets through listening tests in a virtual environment. One of the key findings was that up to five virtual source locations could be reliably identified using stereo panning, regardless of the noise cancellation feature and the spectral content of the sounds. Bone conduction headsets matched the detection accuracy of ANC headsets, despite lower subjective sound quality. Future research will involve visually impaired participants and various distractor sounds to further assess headset performance.

The study by A. M. A. Sabaawi, M. R. Almasaoodi, S. El Gaily, and S. Imre presents the highly constrained quantum genetic algorithm (HCQGA), a novel quantum computing method aimed at solving optimization problems with extremely large and complex search spaces that are impractical for current classical or quantum processors. The HCQGA was tested on maximizing energy efficiency in an uplink multicell massive MIMO system, achieving faster convergence to the optimal solution than classical algorithms.

In their paper, Balázs Ádám Toldi and Imre Kocsis propose a novel blockchain-based method for preserving the confidentiality of collaborative business processes orchestrated by smart contracts, which safeguards sensitive information by only storing encrypted and hashed process states and using zero-knowledge proofs to validate participant actions. This approach, which focuses on a practical subset of BPMN models, ensures secure message-passing between participants and includes an open-source prototype that automatically generates essential software components. The presented solution is unique and has been functionally validated and assessed for its efficiency in terms of computational resources and associated blockchain transaction costs.

Adrian Pekar et al. investigate the adaptive gradual flow aggregation in network flow metering systems as a means to balance the preservation of detailed flow information with the need for data volume reduction, due to the challenges of scalability in managing growing flow entries. Their paper concludes that this method can optimize the trade-off between resource constraints (CPU and memory usage) and the integrity of flow information, with one specific buffer (B2) offering the best balance with a minimal loss of information deemed acceptable for enhanced resource efficiency.

With this overview, let us see all the Infocommunications Journal papers in the 3nd issue of 2023.

**Pal Varga** is the Head of Department of Telecommunications and Media Informatics at the Budapest University of Technology and Economics. His main research interests include communication systems, Cyber-Physical Systems and Industrial Internet of Things, network traffic analysis, end-to-end QoS and SLA issues – for which he is keen to apply hardware acceleration and artificial intelligence, machine learning techniques as well. Besides being a member of HTE, he is a senior member of IEEE, where he is active both in the IEEE ComSoc (Communication Society) and IEEE IES (Industrial Electronics Society) communities. He is Editorial Board member in many journals, and the Editor-in-Chief of the Infocommunications Journal.

# Analysis of the WRC-23 Agenda Item Concerning the Future Use of the 470-694 MHz Band in Europe

Hussein Taha, Péter Vári, and Szilvia Nagy

*Abstract*—The upcoming World Radiocommunication Conference 2023 (WRC-23), according to the preliminary agenda item 1.5, will decide the long-term use of the 470-694 MHz band in International Telecommunication Union Region 1 and consider regulatory actions for the rest of the decade. This article attempts to inform the debate on the future use of spectrum in the 470- 694 MHz band and the status of broadcasting services at the European level and in many individual member states. This article highlights the role WRCs played by allocation of the previous digital dividend bands for mobile services. It provides an overview of the most recent trends and developments of existing services in the sub-700 MHz band. The article also explores and analyzes options for future use of the 470-694 MHz band in Europe and discusses the benefits and implications of adopting each option. Based on this analysis, it is recommended to introduce a secondary allocation to mobile service in the downlink-only mode alongside traditional broadcasting services in the sub-700 MHz band at the upcoming WRC-23.

*Index Terms*—WRC-23, Digital Dividends, 470-694 MHz, sub-700 MHz band, IMT, BBPPDR, DTT, Broadcasting, PMSE

## I. INTRODUCTION

The International Telecommunication Union (ITU) holds World Radiocommunication Conferences (WRCs) every three to four years with the main objective of reviewing and updating the Radio Regulations based on the preliminary agenda items set by the ITU Council [1]. Subsequent to the WRCs, Radio Regulations are issued to allow new radiocommunication systems and applications access to the radio spectrum while safeguarding the operation of existing radiocommunication services.

This section mainly reviews the historical approach of the ITU at WRCs to regulate the Ultra-High Frequency (UHF) spectrum in ITU Region 1. Figure 1 shows how the UHF spectrum was regulated between broadcasting and mobile services following WRCs.

Hussein Taha is with the Doctoral School of Multidisciplinary Engineering Sciences, Széchenyi István University, Győr, Hungary (e-mail: taha.hussein@sze.hu, hussein.taha91@gmail.com)

Péter Vári is with the Department of Telecommunications, Széchenyi István University, Győr, Hungary (e-mail: varip@sze.hu);
He is also with the National Media and Infocommunications Authority, Budapest, Hungary (e-mail: vari.peter@nmhh.hu).

Szilvia Nagy is with the Department of Telecommunications, Széchenyi István University, Győr, Hungary (e-mail: nagysz@sze.hu).

Fig. 1. Spectrum regulation among services operating in the UHF band in ITU Region 1.

In 2006, based on the Geneva agreement (GE-06), analog television broadcasting was switched off and transitioned to digital television broadcasting in Europe, the Middle East, and Africa [2]. The transition period started on 17 June 2006 and ended on 17 June 2015. This transition provided a significant improvement in the quality of TV broadcasting service and allowed some radio spectrum to be freed up for other purposes, known as the "digital dividends".

The ITU allocated these digital dividends also to mobile services (except aeronautical mobile) on a co-primary basis with broadcasting in two phases. In 2007 following WRC-07, the first digital dividend known as the 800 MHz band, the spectrum from 790 MHz to 862 MHz, was allocated for mobile service [3]. During WRC-12/15, the decision was made to allocate the second digital dividend known as the 700 MHz band, the spectrum from 694 MHz to 790 MHz to mobile service from 2015 [4]-[6].

As decided by WRC-15 Resolution 235, the upcoming WRC-23 will review spectrum usage and study the spectrum requirements of existing services in the 470-960 MHz frequency band in ITU Region 1. This will specifically focus on the spectrum needs of broadcasting and mobile services (excluding aeronautical mobile services), with consideration given to relevant ITU Radiocommunication Sector (ITU-R) studies, recommendations, and reports [5], [6].

Based on the recommendations of the latest WRC-19, there is interest in allocating an additional spectrum for International Mobile Telecommunications (IMT) and Broadband Public Protection and Disaster Relief (BBPPDR) services in a band below 694 MHz [7]. Besides, agenda item 1.5 of the upcoming WRC-23 calls for an assessment of the spectrum needs for

existing services in the sub-700 MHz band, which spans from 470 to 694 MHz, in ITU Region 1 and explores the possibility of granting a new IMT allocation in all or parts of the band [8].

Stakeholder viewpoints varied on potential future scenarios for the best use of frequency band 470-694 MHz in Europe. Historically, this band has been allocated and used primarily by Digital Terrestrial Television (DTT), with Program-Making and Special Events (PMSE) services on a secondary basis (plus a small allocation for use by radio astronomy in some countries). DTT and PMSE provide common values to all European citizens. DTT efficiently delivers linear TV services to huge audiences, even for free. It provides near-universal reach, is dependable in emergencies and crises, and enables broad, targeted information. PMSE offers an essential service that connects people through the digital economy.

On the other hand, WRC-23 is a chance to allocate a new digital dividend band for IMT (such as the 600 MHz band) that could aid in boosting the accessibility of 5G, guarantee future growth and innovation, address the digital divide, and can facilitate the introduction of 6G. The WRC is a frame regulation, as the European Conference of Postal and Telecommunications Administrations (CEPT) and/or the European Union (EU) consider technical harmonization measures essential for the introduction of a service other than broadcasting [9]. Additionally, as a result of the Lamy report, the UHF decision provided long-term investment predictability and stimulated innovation by safeguarding, under Article 4, the availability of the 470-694 MHz band for DTT and PMSE services until at least 2030 in all EU member states [10].

Consequently, the future strategy for the sub-700 MHz band should be flexible enough to respond to evolving both broadcasting and mobile market realities while considering technological and consumer behavior developments. In this regard, the article explores and analyzes potential use options for the future 470-694 MHz spectrum in Europe. The benefits and challenges of each alternative are discussed concerning the relevant studies conducted by stakeholders. Based on this analysis, the article outlines the rationale for adopting each option at the upcoming WRC-23.

The rest of this article is structured as follows. Section II highlights the benefits of using digital dividend bands for IMT and BBPPDR. Section III reviews the current situation and future development of DTT, 5G Broadcasting, and PMSE in Europe. Section IV explores and analyzes the potential future uses of the sub-700 MHz band in Europe. Section V concludes this article.

## II.  BENEFITS OF USING DIGITAL DIVIDEND BANDS FOR IMT AND BBPPDR

The term "IMT" (International Mobile Telecommunications) is used to refer to IMT-2000, IMT-Advanced, and IMT-2020 collectively, i.e., the 3G, 4G, and 5G generations of cellular communications systems [11]. IMT systems seek to provide worldwide communication services, regardless of the location, network, or terminal being utilized (global roaming).

IMT systems operating in the digital dividend bands differ from earlier mobile technologies in several aspects, including

frequency allocations, modulation schemes used in uplink and downlink, and resource block allocations in channel bandwidth [12]. The digital dividend bands are crucial for deploying IMT systems, especially in some developing countries and sparsely populated areas, to assist them in reducing the digital divide with urban areas and achieving digital equality [13]. Digital dividends provide wide coverage in rural areas with new technology. Additionally, the propagation characteristics of lower frequencies improve indoor coverage in built-up areas, making access to communications services more inclusive [13]. Figure 2 compares coverage ranges based on frequencies in the open environment. Notably, base stations operating in the digital dividend bands offer broader coverage than those in any other bands above 1 GHz. However, it is essential to consider that propagation within the frequency range of 700-2600 MHz is influenced by various factors, including diffraction, reflection, scattering, and attenuation. These factors significantly differ from ideal free-space propagation, especially when considering the power-law path loss exponent in different scenarios such as urban, suburban, and outdoor-to-indoor environments.



Fig. 2. Coverage comparison in free space depending on frequency.

In urban areas, the power-law path loss exponent tends to be higher compared to free space propagation. This means that the signal strength decreases more rapidly with increasing distance. The power-law path loss exponent in suburban areas is generally lower than in urban areas but still higher than in free space propagation. Consequently, signal strength attenuates faster than in ideal free space conditions but less severely than in densely populated urban areas. Moreover, when a 5G signal transitions from an outdoor environment to an indoor space, additional challenges arise due to building materials and construction elements. In this scenario, the power-law path loss exponent can be quite high, leading to significant signal attenuation as the signal penetrates walls and windows. Besides, indoor reflections and multipath propagation further complicate signal propagation.

The benefits of broadband spectrum are passed on to customers in a competitive market through lower prices, which can increase the take-up and usage of the services.

Global demand for mobile data traffic and broadband multimedia capacity is rising simultaneously with data usage growth. This endangers mobile broadband in emerging markets, rural areas, and inside buildings. The upcoming WRC-23 is a chance to investigate mobile allocations and possibly identify additional spectrum for IMT in the sub-700 MHz band.

The digital dividend frequencies are not just limited to commercial IMT systems; Many countries in ITU Region 1 have deployed BBPPDR services in certain parts of these frequencies by the EU Decision 2016/687 on the harmonization of the 700 MHz band [14], [15]. BBPPDR refers to radio applications utilized by national authorities or relevant operators in response to public safety and security concerns, including emergency situations [14], [15]. Digital dividend spectrum meets the PPDR requirements in broadband by supporting higher data rates, higher capacity, and enhanced multimedia capabilities.

## III. CURRENT SITUATION AND FUTURE DEVELOPMENT OF DTT, 5G BROADCASTING, AND PMSE IN EUROPE

### A. Digital Terrestrial Television

The phrase "linear television" refers to any television programs transmitted (from one to many) following a predetermined schedule to be received by households at certain times on specific channels [16]. Live television content is transmitted in real-time as it is produced, albeit not always linearly. Thus, linear television is essential for both live and non-live content [16]. Linear television content is delivered by broadcasting platforms (DTT, satellite, cable, IPTV). DTT is the most appropriate and viable way to deliver linear TV content to wide audiences in many developing countries.

Technological innovations and advancements in DTT have improved capacity, spectrum efficiency, and service. Initially, the switch from analog to digital terrestrial television broadcasting has increased spectrum efficiency, content diversity, and improved reception/signal quality. The transition of the transmission standard from Digital Video Broadcasting - Terrestrial (DVB-T) (1997) to DVB-T2 (2009) increased capacity by 50-100% for television services in the same amount of radio spectrum. Digital video encoding and compression technologies have evolved over time from the standard Moving Picture Experts Group 2 (MPEG-2) (1996) to H.264 MPEG-4 (1999, Standard Definition (SD)) to H.265 High Efficiency Video Coding (HEVC) (2013, High definition (HD)) and ultimately H.266 Versatile Video Coding (VVC) (2020); each upgrade uses a smaller proportion of the available capacity for each video stream [17]. The requirement for spectrum for DTT services was reduced by around 25% thanks to Single Frequency Networks (SFNs), which allowed multiple transmitters to simultaneously transmit the same signal over the same frequency channel without utilizing additional frequencies [18]. Any additional changes to DTT platforms to boost efficiency (where practicable) would incur high costs.

These upgrades have permitted the resolution improvement from SD to HD services with the same amount of spectrum. Other improvements in resolution and service quality could be introduced on the DTT platform, such as the introduction of Ultra-High Definition (UHD), 4K, High Frame Rate (HFR), High Dynamic Range (HDR), and interactive broadcast broadband, but they would limit the number of video services offered due to their need for higher bandwidths.

However, the DTT platform alone is no longer sufficient to meet the demands of today's society. An examination of the market share of linear TV content delivery platforms in the EU shows that delivery means are heterogeneous across European TV markets [17]-[20]. Figure 3 shows the share of households by platform by member state [17], [19].



Fig. 3. Share of households by platforms, 2021 [17], [19].

The combined data for the EU27 shows that the shares are almost evenly split across the linear TV platforms. However, a predominant platform emerges for several member states. DTT clearly dominates in Greece (71%), Spain (64%), and Italy (48%). Instead, the cable platform dominates in Finland (66%). Satellite is the dominant platform in Slovakia (57%). Finally, IPTV is the most used platform in France (62%).

On the other hand, the future outlook indicates a decrease in linear TV viewing time, while non-linear viewing will continue to dominate daily TV viewing habits [20], [21]. Figure 4 shows viewing time across platforms by the country for the years 2020 and 2021 for some selected EU member states, the United States and Australia, according to new research from OMDIA [21].



Fig. 4. Viewing time across platforms by country [21].

A clear growth is observed in daily viewing habits for the online long-form (Video on Demand (VoD), Over-the-Top (OTT) services like Netflix, Prime Video, Disney+, etc.) and social media video viewing [21].

Indeed, on a platform-by-platform level, pre-scheduled linear TV content still has an inherent value for audiences in the markets covered. Thus, VoD and OTT services complement linear viewing rather than being a replacement. Consequently, linear and non-linear TV will coexist for the foreseeable future.

However, for linear TV, particularly through DTT, to remain a part of the landscape in the ever-evolving TV market for many years to come, new technologies, such as standalone 5G broadcasting and collaborative content creation for television over various media must be developed to deliver more immersive content to the audience across multiple devices.

*B. 5G Broadcasting*

5G Broadcasting uses 5G technology to deliver linear broadcast content directly to multiple users on mobile devices (such as smartphones, tablets, and cars) without Wi-Fi or the Internet [22], [23]. 5G Broadcasting is seen as a promising supplement to DVB-T2, not a substitute for stationary reception. 5G Broadcasting can offer a similar capacity to DTT and transmit both linear TV and radio programs, assuming good reception conditions. Importantly, standalone 5G Broadcasting can utilize the sub-700 MHz band because it is intended to operate alongside DTT without interference [23], [24].

However, several challenges need to be tackled before 5G Broadcasting can be widely adopted. These issues include the spectrum that will be used (current DTT spectrum or another allocated spectrum), the networks that will be used (existing DTT networks or mobile networks), the requirements for free broadcast, and business models that can support 5G broadcasting. Additionally, since 5G broadcasting combines the broadcasting and telecom industries, creating a centralized approach between them is essential [23], [24].

Trials to enhance the 5G broadcasting standard are now underway in Europe. The first phase of the 5G broadcasting trial in Vienna ended, which ran from Q4/2019 to Q2/2021 and compared the 5G broadcasting with DTT [25]. They are currently in the process of the second phase, which will extend until Q3/2023 and aims to investigate the applications and further develop the 5G broadcasting ecosystem [26]. This trial reached the following conclusions [25], [26]. 5G broadcasting significantly extends the reach of terrestrial broadcasting as enabled devices can be used for a portable outdoor reception. 5G broadcasting can achieve comparable performance compared to DVB-T2, with potential enhancements in upcoming 3GPP versions. 5G broadcasting networks supplement existing DVB-T2 networks for fixed and portable indoor reception. Consequently, 5G broadcasting is suitable to coexist with DTT in the sub 700 MHz band and enables innovation in terrestrial broadcasting, although the business models that could support 5G broadcasting are still unclear.

Given the concerns about the business models that could support 5G broadcasting, there is an opportunity to benefit from the factors essential for defining the business model and pricing of future industrial mobile networks, as outlined in [27]. This contribution delves into the classification of networks based on customer requirements, constraints, and motivations. It also presents illustrative use cases for each scenario, including the

business advantages of cloud solutions, challenges related to frequency allocation, the potential of network slicing, and the importance of energy-efficient networks [27].

*C. Program-Making and Special Events*

Today, alongside DTT, the sub-700 MHz band in the UHF spectrum is widely used throughout Europe by PMSE equipment, such as wireless microphones, in-ear systems, talkback systems, camera control systems, audio/video links, and so forth [17], [28]. PMSE provides content production services to broadcasters, professional content producers, and various organizations in the community.

The demand for content created by PMSE is increasing steadily, driven by both the current audiences and the expanded worldwide audience realized by new delivery platforms. Moreover, the PMSE industry is characterized by innovation, as seen by the adoption of IP and cloud-based workflows, the introduction of digital audio technology, the usage of 5G for some PMSE applications, and the growing complexity of production [28]. Thus, it is difficult to quantify and average the required spectrum given the wide range of activities using PMSE equipm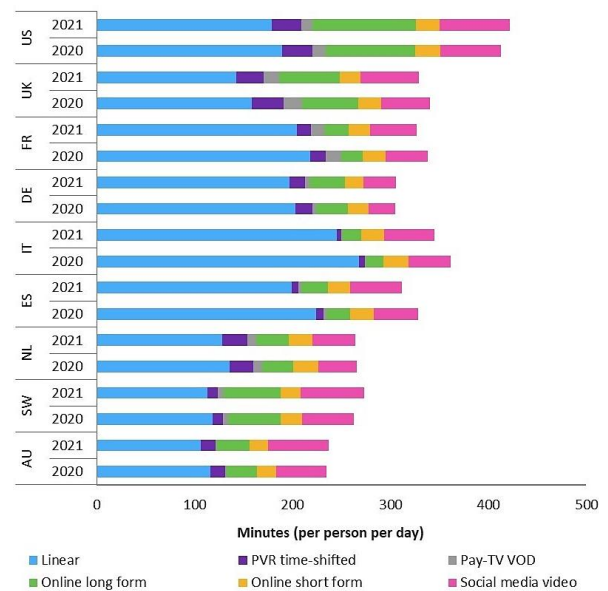ent and the increasing demand for complicated and advanced productions. Additionally, the need for spectrum will significantly increase during special events such as championships, elections, and other large events.

The UHF band has the best conditions to meet the needs and demands of PMSE in terms of high spectrum efficiency, ultra-low latency, high transmission reliability, and high audio quality. Future use of this band for PMSE services will rely on decisions taken at the upcoming WRC-23.

## IV. POSSIBLE OPTIONS FOR FUTURE USE OF THE 470-694 MHZ BAND IN EUROPE

Decisions taken at WRC-23 will affect future EU policies on using the 470-694 MHz band in Europe. Any EU action should adhere to the ITU Radio Regulations and the EU's policy objectives in order to promote the development of a single digital market supported by dependable high-speed networks and encourage the effective management and use of radio spectrum. This section considers four options for the future use of the 470-694 MHz band in Europe. Stakeholder viewpoints for each option are investigated.

*A. Option 1 – no change in regulation sub-700 MHz band.*

In this scenario, Europe will take no change on the frequency band 470-694 MHz. Mainly broadcasting services will continue to use the whole 224 MHz spectrum segment. No spectrum is available in the sub-700 MHz band for mobile services within the framework of EU harmonization. Figure 5 shows the sub-700 MHz band distribution according to option 1.

Fig. 5. Option 1 – no change in regulation the sub-700 MHz band.

European Broadcasting Union (EBU), and Broadcast Network Europe (BNE), support the "no change" option, justifying their viewpoint as follows [29]-[31]. Currently, DTT and PMSE in the sub-700 MHz band serve the community's public and commercial needs and may continue to be required to deliver public service media content well in many EU countries for many years.

Under the EU Decision 899/2017, the European Parliament, Council, and Commission have guaranteed access to the band 470-694 MHz for terrestrial broadcasting services until at least 2030 [2], [32]. Full access to the 470-694 MHz band for DTT with PMSE will maximize the economic benefits to the EU and ensure efficient spectrum use. Additionally, this will promote innovation and help broadcast networks use energy-efficient and environmentally friendly broadcasting. One of the key results of the ITU report BT.2302-1 is that most countries express a need for the whole 224 MHz spectrum segment for broadcasting [33]. Only Slovenia and Finland in Europe have declared a desire for less than 224 MHz.

On the other hand, the analysis of responses to ITU-R circular letter 6/LCCE/104 at task group TG 6/1 shows that the current spectrum available for broadcasting in the sub-700 MHz band is sufficient and essential for the broadcasting services, and most of the responding administrations have indicated that they plan to continue using this band for broadcasting services [34].

However, IMT is currently in a situation with increased demand for more low-band spectrum to provide cost-effective coverage solutions. In this "no-change" scenario, further investment in network infrastructure might address the current coverage inadequacies rather than additional spectrum in the UHF band.

Mobile network capacity can be increased through network densification and using previously available digital dividend bands [35]. The development of fiber networks complemented by Wi-Fi connectivity is also expected to meet a significant portion of the goal of digital connectivity in the digital decade [35]. Furthermore, reorganizing the sub-1 GHz mobile bands and phasing out outdated mobile operator technologies like 2G from 900 MHz will permit big contiguous frequency blocks, allowing their effective capacity to be almost trebled [36].

*B. Option 2 – co-primary allocation to mobile service in the sub-700 MHz band.*

In this scenario, the CEPT and/or EU would consider harmonized technical solutions on a new digital dividend band to the mobile service, for example, the "600 MHz band" (606-694 MHz), while the spectrum (470-606 MHz) would still be available for DTT and PMSE. Figure 6 shows the frequency allocation of the sub-700 MHz band according to option 2.

There is a consensus that mobile traffic is growing worldwide, and therefore more spectrum is needed to increase the capacity of mobile networks. Given the historical allocation of prior digital dividend bands for mobile purposes during previous WRCs, could raise expectations that the 606-694 MHz spectrum may be taken away from broadcasting services if a co-primary allocation is agreed upon at WRC-23.



Fig. 6. Option 2 – co-primary allocation to mobile service in the sub-700 MHz band.

The Groupe Special Mobile Association (GSMA) supports allocating additional spectrum for IMT, arguing that access to the 600 MHz band has several advantages [13], [37]. The additional spectrum in the 600 MHz band may reduce the cost of extending coverage to rural residents by 33% and improve rural and deep-indoor broadband speeds by 30-50% [37]. Utilizing the 600 MHz frequency for specific broadband services like BBPPDR could also benefit society.

However, introducing a co-primary allocation for 470-694 MHz to IMT at WRC-23 would also lead to negative impacts. The loss of any further spectrum is likely to undermine the public and commercial value that DTT and PMSE services already provide in the sub-700 MHz band, as well as inhibit the innovation and investments necessary to maintain and grow this value in the future [29], [30]. Indeed, according to the report on spectrum inventory by the European Commission, the demand for future spectrum consumption was at the same level for both the broadcasting and mobile industries in the short, medium, and long term [38].

On the other hand, enabling the fragmented approach to the sub-700 MHz band would require more time and effort on a large scale throughout Europe in terms of planning, preparations, international frequency coordination, technical activities, as well as the cost of migrating DTT transmitters onto new frequencies. Various studies have demonstrated that mobile and broadcasting services cannot operate on the same frequencies in the same or nearby areas without causing mutual harmful interference [39], [40]. In practice, if a country used the 600 MHz band for mobile purposes before its neighbors, it could cause interference to DTT users in neighboring countries in areas up to tens of kilometers from the border, representing about 13% of the EU population [41].

This issue would be particularly important for EU members that border non-EU countries, like Hungary. They have to switch off the broadcasting services in the 600 MHz band, but they cannot launch IMT until the necessary cross-border coordination agreements with their neighbors have been concluded. GE-06 agreement and ITU-R report BT.2383 include the final acts for planning and protecting DTT in cross-border territories [2], [42].

Thus, EU action will have to include effective management of interference with DTT outside the EU borders. Interference can be reduced by geographically separating mobile and broadcasting services. However, the required separation distances can sometimes be hundreds of kilometers. Compatibility studies of IMT and broadcasting services in the prior digital dividend bands can benefit administrations in developing bilateral agreements for the new band [39]-[45].

*C. Option 3 – secondary allocation to mobile service in the downlink-only mode in the sub-700 MHz band.*

In this scenario, Europe would consider a flexible use that allows the coexistence of mobile broadband services in the downlink-only mode (i.e., unidirectional) alongside traditional broadcasting services in the sub-700 MHz band, where there is no or decreasing demand for DTT at the national level. Figure 7 shows the frequency allocation of the sub-700 MHz band according to option 3.

Fig. 7. Option 3 – flexibility approach to the sub-700 MHz band.

In practice, each country protects DTT's sub-700 MHz band spectrum access to the extent deemed necessary, while channels not used for DTT might become available for downlink-only services or applications served by mobile broadband technologies, depending on national conditions. To provide PMSE services, the "white spaces" in the sub-700 MHz band will be used.

Several broadcasting and mobile manufacturers support this flexibility option, justifying their viewpoint as follows [10], [46]. According to national demand, a combination of wireless broadband services in downlink-only mode and/or DTT services will be offered in the sub-700 MHz band to maximize positive economic benefits and minimize negative social repercussions.

By adopting this option, mobile operators can enhance their downlink capacity to satisfy the growing demand for mobile traffic where spectrum is available while also ensuring the availability of the necessary spectrum for the future development of DTT in Europe. Relevant statistics have indicated that traffic requires significantly more downlink capacity than uplink capacity, owing primarily to increased video and app-based mobile use [47].

Additionally, allowing only downlink services will limit the fragmentation of UHF spectrum utilization and sufficiently mitigate interference between IMT and DTT services. Uplink services won't be permitted since they would entail significant service restrictions due to the interference mitigation with broadcasting services needed to comply with the GE-06 agreement and the ITU Radio Regulations, particularly at the EU outer borders.

However, member states that allow IMT-only downlink services will have to guarantee cross-border coordination with neighboring countries where DTT is operational. Studies should be conducted to determine the technical coexistence conditions of DTT and IMT in the downlink-only mode in the sub-700 MHz band.

On the other side, when option 3 is adopted, the densification of DTT networks will make it more difficult for PMSE users to access "white spaces".

*D. Option 4 – allocation to the mobile service in the whole sub-700 MHz band.*

In this option, the regional coordination will define a common roadmap for clearing DTT from the sub-700 MHz band and allocating the full 224 MHz spectrum segment to mobile broadband services and possibly other sectorial services such as PMSE and BBPPDR. Figure 8 shows the frequency allocation of the sub-700 MHz band according to option 4.

Fig. 8. Option 4 – allocating the whole sub-700 MHz band for mobile services.

In practice, the 470-694 MHz spectrum will be available to mobile broadband services, and DTT transmission will cease. Thus, in the absence of the DTT platform, there are two possible scenarios for providing linear broadcast television services to all EU households in future.

The first scenario would be based on the migration of DTT viewing to alternative platforms (a mixture of satellite, cable, and IPTV). Aetha Consulting Ltd has considered this alternative scenario and estimated the economic costs and benefits arising for EU citizens over a 15-year period (from 2015 to 2029) [16].

The potential costs associated with migrating DTT viewing to alternative platforms encompass several aspects. Firstly, there are costs related to acquiring client premises equipment and transponder capacity necessary to deliver linear TV content on these alternative platforms. Another significant cost implication lies in the potential increase in expenses for broadcast companies. With fewer platforms competing for viewers after the migration, broadcasters may face higher costs as they strive to retain their audience and remain competitive in the evolving media landscape. In countries like Spain and Italy, where many local channels are provided through the DTT platform, ceasing DTT could result in losing access to these local TV services. In countries like Spain and Italy, where many local channels are provided through the DTT platform, ceasing DTT could result in losing access to these local TV services. Additionally, there may be a loss due to consumers' preference for DTT, as it provides free-to-view services. On the other hand, when the DTT platform is ceased, the spectrum previously allocated for broadcasting may become unavailable for PMSE use. As a result, moving PMSE users to new frequencies could become more expensive and logistically challenging.

Whereas the availability of the 470-694 MHz band for mobile services could bring about several significant benefits. Mobile operators stand to gain cost savings since they won't have to invest in deploying new base station sites to increase capacity, and thus, these savings will be passed on to consumers through lower prices. Moreover, there will be additional savings from no longer needing to maintain and operate the DTT network, including expenses related to power, staff,

equipment maintenance, and future equipment upgrades. Another advantage is the potential for faster adoption of high-speed broadband services, which could effectively address digital inequalities and bridge the digital divide. On the other side, the 470-694 MHz spectrum could also be utilized by other broadband services like BBPPDR.

The economic evaluation of this scenario indicated that the costs associated with migrating DTT viewing to alternative platforms would be many multiples of the benefits that would arise from making the whole sub-700 MHz band available for mobile services [16]. However, this assessment may differ between EU member states depending on the current levels of adoption of existing television platforms.

The second scenario explores and evaluates the possibility of delivering television channels currently provided through DTT via Mobile Broadband (MBB) networks. This scenario is investigated as a case study in Finland [48], [49]. The objective is to assess the required investment to continue providing TV channels via MBB networks in Finland by 2030. To employ this scenario, public service broadcasting must be accessible to everyone through the open Internet, technically functional, and be cost-free at the point of use.

Unicast delivery of linear TV content to all viewers requires a significant investment in rural unicast capacity. The existing network grids would need to be densified to provide uniform 5G coverage on higher bands. Alternatively, evolved Multimedia Broadcast Multicast Service (eMBMS) or 5G Xcast would be required to optimize the low-band capacity in mobile broadband networks [48]. However, with 4G and 5G mobile networks, broadcast and multicast functionality can be deployed together with typical 4G/5G unicast. This means the same content can be delivered to numerous users with one transmission, effectively saving spectrum capacity when multiple viewers watch the same content simultaneously [49].

The main challenge for this scenario is securing wireless broadband capacity during peak viewing times in sparsely populated rural areas. On the other hand, PMSE technologies will need to adapt to changes in the available spectrum.

## V. CONCLUSION

Future EU policy on the 470-694 MHz spectrum after 2030 is expected to consider decisions taken at the upcoming WRC-23. This article provided an overview of the most recent trends and advancements concerning the future uses of the sub-700 MHz spectrum in Europe.

Analysis of the current situation and future development of broadcasting services currently operating in the 470-694 MHz band in the EU showed that while linear TV watching, mainly through DTT, continues to be popular with EU many member states, the demand for VoD services and OTT subscriptions are increasing to adapt to user behavior.

Consequently, the TV future will result in higher quality, improved sound, interactive TV, VoD and OTT offerings that complement DTT (DTT for live TV viewing, Hybrid TV through VoD and OTT for non-linear TV), as well as 5G broadcasting that complements DTT (DTT at home, 5G broadcasting on the move).

Analysis of options for future use of the 470-694 MHz band showed that all future strategies should be able to adjust to the reality of broadcasting and the growing mobile market while considering technology advancements and consumer behavior.

The continued use of the whole sub-700 MHz band for broadcasting services, as envisaged in option 1, will ensure the spectrum needed to support further growth and innovation in broadcasting services. Nevertheless, this will delay the deployment of mobile broadband networks, particularly in border rural areas where coverage will be poor or nonexistent. Additionally, border urban areas may experience capacity limitations for high-speed broadband due to a lack of spectrum.

Allocating new digital dividend bands to IMT and BBPPDR systems in the sub-700 MHz band would help lower broadband prices, increase accessibility to communications services, and overcome the digital divide. However, enabling the fragmented approach to the 470-694 MHz band when allocating the new 600 MHz band to the mobile service on a primary basis, as envisaged in option 2, would undermine DTT and PMSE services and necessitate the EU countries adopt a national roadmap and reach all necessary cross-border coordination agreements while adhering to the transition deadline.

Whereas launching the mobile service in the downlink-only mode in the sub-700 MHz band on a secondary basis, as envisaged in option 3, would protect the interest of both broadcasting and mobile services in the 470-694 MHz band and allow the development of an innovation ecosystem for both. Banning mobile services in uplink mode would sufficiently mitigate the mutual interference between IMT and DTT systems. Therefore, it is recommended that this flexibility option be adopted at WRC-23. However, compatibility studies are still required to determine the technical coexistence requirements of DTT, PMSE, and IMT in the downlink-only mode in the sub-700 MHz band at the EU level.

On the other hand, the full-scale reallocating, as envisaged in option 4, is not currently considered suitable because DTT and PMSE continue to play a significant role in EU many member states.

## REFERENCES

[1] International Telecommunication Union (ITU), Geneva, Switzerland. [Online]. Available: https://www.itu.int/

[2] ITU, "Final Acts of the Regional Radiocommunication Conference for planning of the digital terrestrial broadcasting service in parts of Regions 1 and 3, in the frequency bands 174-230 MHz and 470-862 MHz (RRC-06)," Geneva, Switzerland, 15 May-16 June 2006 [Online]. Available: https://www.itu.int/pub/R-ACT-RRC.14-2006/en

[3] ITU, "The Final Acts of the World Radiocommunication Conference 2007 (WRC-07)," Geneva, Switzerland, 2007 [Online]. Available: http://handle.itu.int/11.1002/pub/802313bd-en

[4] ITU, "The Final Acts of the World Radiocommunication Conference 2012 (WRC-12)," Geneva, Switzerland, 2012 [Online]. Available: http://handle.itu.int/11.1002/pub/805627a4-en

[5] ITU, "The Final Acts of the World Radiocommunication Conference 2015 (WRC-15)," Geneva, Switzerland, 2015 [Online]. Available: http://handle.itu.int/11.1002/pub/80d4e1c0-en

[6] The European Parliament and the council of the European Union, "Decision (EU) 2017/899 of the European Parliament and of the Council of 17 May 2017 on the use of the 470-790 MHz frequency band in the Union," Strasbourg, France, 2017. pp. 1–7 [Online]. Available: http://data.europa.eu/eli/dec/2017/899/oj

[7] ITU, "The Final Acts of the World Radiocommunication Conference 2019 (WRC-19)," Sharm El-Sheikh, Egypt, 2019 [Online]. Available: http://handle.itu.int/11.1002/pub/813b5921-en

[8] ITU, World Radiocommunication Conference 2023 (WRC-23), Agenda items and status of studies, Dubai, United Arab Emirates, 20 November to 15 December 2023 [Online]. Available: https://www.itu.int/wrc-23/

[9] Latest CEPT Briefs and ECPs developed by CPG in preparation of WRC-23, last updated: 21 February 2023 [Online]. Available: https://www.cept.org/ecc/groups/ecc/cpg/page/cept-briefs-and-ecps-for-wrc-23/

[10] P. Lamy, "Results of the work of the high-level group on the future use of the UHF band (470-790 MHZ)," Report to the European Commission 1, 2014, pp. 1–34 [Online]. Available: https://ec.europa.eu/newsroom/dae/document.cfm?doc_id=6721

[11] ITU-R, Naming for International Mobile Telecommunications, Resolution 56-2 (2015) [Online]. Available: https://www.itu.int/pub/R-RES-R.56-2-2015

[12] H. Taha, P. Vári, and S. Nagy, "On the Challenges of Mutual Interference between Cable Television Networks and Mobile Fixed Communication Networks in the Digital Dividend Bands," Infocommunications Journal, Vol. XIV, No 3, September 2022, pp. 63–71. DOI: 10.36244/ICJ.2022.3.8

[13] Coleago Consulting Ltd and GSMA, Low-band spectrum for 5G report, the need for sub-1 GHz spectrum to deliver the vision of 5G, 9 May 2022, pp. 1–32 [Online]. Available: https://www.gsma.com/spectrum/wp-content/uploads/2022/07/Low-Band-Spectrum-for-5G.pdf

[14] European Commission, Commission Implementing Decision (EU) 2016/687 on the harmonisation of the 694-790 MHz frequency band for terrestrial systems capable of providing wireless broadband electronic communications services and for flexible national use in the Union, Brussels, 28 April 2016, pp. 1–12 [Online]. Available: http://data.europa.eu/eli/dec_impl/2016/687/oj

[15] Resolution 646 (Rev.WRC-19), Public protection and disaster relief, The World Radiocommunication Conference (WRC-2019), Sharm El-Sheikh, Egypt, 2019, pp. 443–448 [Online]. Available: https://www.itu.int/dms_pub/itu-r/oth/0C/0A/R0C0A00000F00133PDFE.pdf

[16] Aetha Consulting Limited, "Report on the future use of the 470-694 MHz band in the EU," Cambridge, United Kingdom, 31 October 2014, pp. 1–89 [Online]. Available: https://aethaconsulting.com/wp-content/uploads/2021/05/Future-use-of-the-470-694MHz-band-in-the-EU-October-2014.pdf

[17] European Commission, Directorate-General for Communications Networks, Content and Technology, Study on the use of the sub-700 MHz band (470-694 MHz): final report, Publications Office of the European Union, 2022, pp. 1–273 [Online]. Available: https://data.europa.eu/doi/10.2759/94757

[18] EBU Technical Report 016, Benefits and limitations of single frequency network (SFN) for DTT, October 2012, pp. 1–9 [Online]. Available: https://tech.ebu.ch/docs/techreports/tr016.pdf

[19] European Audiovisual Observatory, 2021, [Online]. Available: https://www.obs.coe.int/en/web/observatoire/

[20] L. Chiariglione and C. A. Szabó, "Multimedia Communications: Technologies, Services, Perspectives, Part II. Applications, Services and Future Directions," Infocommunications Journal, Vol. VI, No 3, September 2014, pp. 51–59., [Online]. Available: https://www.infocommunications.hu/documents/169298/831299/InfocomJ_2014_3_8_Chiariglione.pdf

[21] R. Moyser and A. Thomas, Cross-Platform Television Viewing Time Report – 2022, OMDIA, London, UK, 22 July 2022, [Online]. Available: https://omdia.tech.informa.com/OM025343/Cross-Platform-Television-Viewing-Time-Report--2022

[22] ETSI and EBU, Technical Specifications 103 720 V1.1.1, 5G Broadcast System for linear TV and radio services; LTE-based 5G terrestrial broadcast system, 2020, pp. 1–44 [Online]. Available: https://www.etsi.org/deliver/etsi_ts/103700_103799/103720/01.01.01_60/ts_103720v010101p.pdf

[23] EBU Technical Report 063, 5G Broadcast Network Planning and Evaluation, Geneva, August 2021, pp. 1–93 [Online]. Available: https://tech.ebu.ch/publications/tr063

[24] EBU Technical Report 064, Compatibility between 5G Broadcast & other DTT systems in the sub-700 MHz band, Geneva, August 2021, pp. 1–47 [Online]. Available: https://tech.ebu.ch/publications/tr064

[25] CEPT, Input Paper, 5G Broadcast Vienna Field Trial, Vienna, Austria 2021, pp. 1–11 [Online]. Available: https://www.ors.at/fileadmin/user_upload/ors/5G_Broadcast/PTD_21_012_Technical_implementation_status_of_5G_Broadcast_-_Vienna_Field_Trial_.pdf

[26] ORS group, Austrian Broadcasting Services (ORS), 5G Broadcast, Vienna, Austria, [Online]. Available: https://www.ors.at/en/5g-broadcast/

[27] G. Soós, D. Ficzere, T. Seres, S. Veress and I. Németh, "Business opportunities and evaluation of non-public 5G cellular networks – A Survey", Infocommunications Journal, Vol. XII, No 3, July 2020, pp. 31–38. DOI: 10.36244/ICJ.2020.3.5

[28] ECC Report 323, Spectrum use and future requirements for PMSE, Denmark, 12 February 2021, pp. 1–54 [Online]. Available: https://docdb.cept.org/download/3470

[29] European Broadcasting Union (EBU), white paper, "'No Change' at WRC-23 maximizes public value and innovation in the UHF band," Geneva, November 2021, pp. 1–6 [Online]. Available: https://tech.ebu.ch/publications/white-paper-20.11.2021

[30] Broadcast Networks Europe (BNE), "'No change' to the UHF band at WRC-23 enables the ongoing success of essential broadcasting services," Brussels, May 2022, pp. 1–6 [Online]. Available: https://broadcast-networks.eu/wp-content/uploads/No-change-on-WRC23-AI1.5-BNE-position-paper.pdf

[31] European Conference of Postal and Telecommunications (CEPT), "Draft CEPT Brief on WRC-23 agenda item 1.5," Yverdon-les-Bains, Switzerland, 07–11 November 2022, pp. 1–57 [Online]. Available: https://cept.org/Documents/cpg/74230/cpg-22-042-annex-iv-05_draft-cept-brief-on-wrc-23-agenda-item-15

[32] European Parliament and the Council of the European Union, "Decision (eu) 2017/899 of the European Parliament and the Council of 17 May 2017 on the use of the 470-790 MHz frequency band in the Union," Strasbourg, France, 17 May 2017, pp. 1–7 [Online]. Available: https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32017D0899

[33] ITU-R report BT.2302-1, "Spectrum requirements for terrestrial television broadcasting in the UHF frequency band in Region 1 and the Islamic Republic of Iran," Geneva, Switzerland, 2021, pp. 1–54 [Online]. Available: https://www.itu.int/pub/R-REP-BT.2302-1-2021

[34] Radiocommunication Bureau (BR), Circular letter 6/LCCE/104, Questionnaire on spectrum use and spectrum needs for terrestrial television broadcasting in the UHF frequency band in connection with WRC-23 agenda item 1.5, Geneva, Switzerland, 27 February 2020 [Online]. Available: https://www.itu.int/dms_pub/itu-r/md/00/sg06/cir/R00-SG06-CIR-0104!!PDF-E.pdf

[35] European Commission, Digital Economy and Society Index (DESI), 2020, page 44 and page 15 [Online]. Available: https://ec.europa.eu/newsroom/dae/document.cfm?doc_id=67086

[36] GSMA, Legacy mobile network rationalization; Experiences of 2G and 3G migrations in Asia-Pacific, May 2020, pp. 1–62 [Online]. Available: https://www.gsma.com/spectrum/wp-content/uploads/2020/06/Legacy-mobile-network-rationalisation.pdf

[37] Coleago Consulting Ltd and the GSMA, Vision 2030: Low-Band Spectrum for 5G, June 2022, pp.1–20 [Online]. Available: https://www.gsma.com/spectrum/wp-content/uploads/2022/07/5G-Low-Band-Spectrum-1.pdf

[38] European Commission, "Report from the Commission to the European Parliament and the Council on the Radio Spectrum Inventory," COM (2014) 536 final, Brussels, Belgium, 1 September 2014, pp. 1–14 [Online]. Available: https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex%3A52014DC0536

[39] ITU-R report BT.2337, "Sharing and compatibility studies between digital terrestrial television broadcasting and terrestrial mobile broadband applications, including IMT, in the frequency band 470-694/698 MHz," Geneva, 2017, [Online]. Available: https://www.itu.int/pub/R-REP-BT.2337

[40] ITU-R report BT.2301, "National field reports on the introduction of IMT in the bands with co-primary allocation to the broadcasting and the mobile services," Geneva, Switzerland, 2022, [Online]. Available: https://www.itu.int/pub/R-REP-BT.2301

[41] ESPON & University of Geneva, "European Perspective on Specific Types of Territories," 2013, pp. 1–140 [Online]. Available: https://www.espon.eu/sites/default/files/attachments/GEOSPECS_Final_Report_v8___revised_version.pdf

[42] ITU-R report BT.2383, "Typical frequency sharing characteristics for digital terrestrial television broadcasting systems in the frequency band 470-862 MHz," Geneva, Switzerland, 2022 [Online]. Available: https://www.itu.int/pub/R-REP-BT.2383

[43] ITU-R report BT.2338, "Services ancillary to broadcasting/services ancillary to programme making spectrum use in Region 1 and the implication of a co-primary allocation for the mobile service in the frequency band 694-790 MHz," Geneva, Switzerland, 2014 [Online]. Available: https://www.itu.int/pub/R-REP-BT.2338

[44] ITU-R report BT.2339, "Co-channel sharing and compatibility studies between digital terrestrial television broadcasting and international mobile telecommunication in the frequency band 694-790 MHz in the GE06 planning area," Geneva, Switzerland, 2014 [Online]. Available: https://www.itu.int/pub/R-REP-BT.2339

[45] H. Taha, P. Vári, and S. Nagy, "Survey on Coexistence of Terrestrial Television Systems and Mobile Fixed Communications Networks in Digital Dividend Bands," *Proceedings of the Technical University of Sofia (TU-Sofia)*, Bulgaria, 2022, pp. 29–36
**DOI**: 10.47978/TUS.2022.72.01.005

[46] RSPG Secretariat, "RSPG Opinion on a long-term strategy on the future use of the UHF band (470-790 MHz) in the European Union," Brussels, Belgium, 19 February 2015, pp. 1–34 [Online]. Available: https://circabc.europa.eu/d/a/workspace/SpacesStore/c88af26c-0f24-4431-8647-56deaa917307/RSPG15-595_final-RSPG_opinion_UHF.pdf

[47] Data Report, We Are Social and Hoot suite, "Distribution of global monthly mobile data volume as of January 2021, by category," Chart, 27 January 2021, Statista, [Online]. Available: https://www.statista.com/statistics/383715/global-mobile-data-traffic-share/

[48] Omnitele Ltd and Yle Finnish Broadcasting Company, the first study, "Assessment for YLE; Linear content in mobile networks; Case: Finnish mobile networks and unicast delivery," Finland, May 2019, pp.1–58 [Online]. Available: https://drive.google.com/file/d/1mMA4ZAeXCMJfRzcbWeuUtcxzt4jd-TV0

[49] Omnitele Ltd and Yle Finnish Broadcasting Company, the second consultant study, "Assessment for Yle; Broadcast Investment Assessment in Mobile Broadband Networks," Finland, December 2019, pp. 1–81 [Online]. Available: https://drive.google.com/file/d/19HvvBxq1fB7D6q3jYSNlv7mdhDHiYKfg/view?usp=sharing

**Hussein Taha** is a PhD candidate at the Doctoral School of Multidisciplinary Engineering Sciences at Széchenyi István University in Hungary. He holds a BSc in Communication and Electronics Engineering from Tishreen University in Syria and MSc in Telecommunication Engineering from the same university in 2019. His areas of interest are Mobile and Wireless Communications, and Broadcasting.

**Péter Vári** is an AssociateProfessor in the Department of Telecommunications at Széchenyi István University, the university of Győr in Hungary. He is now Deputy Director-General for Technical Affairs at the National Media and Infocommunications Authority in Hungary. His general interests span the areas of Radiocommunications, Mobile Services, and Broadcasting.

**Szilvia Nagy** is a Professor in the Department of Telecommunications at Széchenyi István University, the university of Győr in Hungary. Her research interests include Digital Image Processing, Information and Coding Theory, DSP , EMC, Quantum Information Theory, and Semiconductors.

# Picosatellite identification and Doppler estimation using passive radar techniques

Tibor Herman and Levente Dudás

*Abstract*—In this article a novel method for satellite identification is presented for picosatellites. Utilizing apriori knowledge of the transmitted signals the cross-correlation of the received signal and a known transmission is calculated, from the results the Doppler shift is estimated and the satellite's Doppler curve can be matched against the measurements of NORAD. Using this method an omnidirectional antenna can be used instead of a high gain directional antenna, and the gain difference is compensated by the processing gain of the algorithm described. The practicality of the algorithm is demonstrated through the mission of MRC-100 PocketQube.

*Index Terms*—PocketQube, satellite identification, remote sens- ing, signal processing, passive radar

## I. INTRODUCTION

SINCE 2021 the popularity of the picosatellite category called PocketQube has been rising rapidly, with a growing number of them being deployed [1]–[5]. They are usually put to orbit in groups, which means that the small distance between them and their radar cross section makes the individual identification by ground based radar systems challenging. The base unit of such spacecraft is a 5 cm cube, which is close to the minimum detection size of the North American Aerospace Defense Command's (NORAD) radars [6]. Especially at the beginning of a mission the two-line element set (TLE) provided by them can be quite inaccurate, because they are not updated every day and there are several unidentified sets that correspond to the targets close to each other. This, in conclusion degrades the reception quality of the ground stations, because the Doppler estimation relies on these measurements.

Since the available onboard power of a PocketQube is limited, the transmission power is usually in the range of 100 mW or 20 dBm. It produces a received power level on the ground station that has little margin in the link budget, so usually the tracking of such satellites is done using directional antennas, which have to track the movement of the spacecraft. The accuracy of tracking is mainly defined by the available TLE provided by NORAD [8].

The aim of this paper is to provide a means for identifying pico satellites without an expensive tracking and antenna system, using a simple omnidirectional antenna, such as a ground plane or dipole. The gain lost due to the lack of tracking and antenna gain is recovered by the processing gain introduced during signal processing. In the following sections,

T. Herman and L. Dudás are with Budapest University of Technology and Economics, Department of Broadband Infocommunications and Electromagnetic Theory, Faculty of Electrical Engineering and Informatics, Budapest, Hungary. (e-mail: herman.tibor@vik.bme.hu, dudas.levente@vik.bme.hu)

the latest satellite mission of the university will be used to demonstrate the practicality of the algorithm described.

There have been little research on the topic, the most relevant paper in this field detects the reflected signal of GRAVES (French: Grand Réseau Adapté à la Veille Spatiale), which is a French radar surveillance system, transmitting on 143,05 MHz [9]. In this paper the authors were able to detect satellites on low Earth orbit in the range of several hundred kilograms.

In [10], [11] the authors use the forward scattered signals of a broadcast satellite to detect high altitude targets, like a space shuttle and they emphasize the importance of the large radar cross-section (RCS) required for this technique to work. A picosatellite, like a PocketQube has much smaller RCS, so although it seems possible, research has to be done to find out about the feasibility regarding small targets.

To increase the quality of reception we propose an algorithm that measures the Doppler shift of the received signal by cross-correlating the received samples with a predefined pattern that is transmitted periodically by the on board radio. This method is usually used in passive radar applications, where a radio receiver is used to measure the distance and relative velocity of a target without transmitting anything. It relies on ambient signals that are transmitted by high power broadcast stations. The reference signal received by the observer is compared with its Doppler shifted and time delayed copies to produce a correlation peak in the ambiguity function [12], [13].

MRC-100 (Fig. 1) is a 3 unit PocketQube satellite developed by students and lecturers of Technical University of Budapest. It was named after Műegyetemi Rádió Club, honoring its $100^{th}$ anniversary in 2024. It is the latest member of the SMOG satellites, whose task is to monitor human caused electromagnetic pollution on a low Earth orbit [15], [16]. The satellite transmits packets periodically on 436.72 MHz using Gaussian Minimum Shift Keying (GMSK) modulation [8]. Non-coherent synchronization bits found in the packets are used as a reference signal which are compared against the incoming signal, while the incoming signal is delayed and shifted in frequency. The output of the correlation function has a global maximum where the two signals match in the delay-Doppler matrix, assuming that the signal we are looking for is present in the received signal.

## II. SIGNAL DETECTION

To detect a predefined signal in a received noisy signal that is shifted due to the Doppler effect, the 2D cross-correlation function of the two signals have to be calculated while one of the signals is copied with Doppler shift. Its discrete time
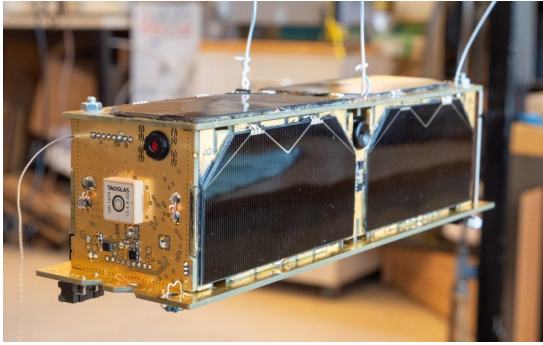
Picosatellite identification and Doppler estimation
using passive radar techniques



Fig. 1. MRC-100 flight model

definition is given by Eq. 1, where $n$ is the sample delay, $f_D$ is the Doppler shift, $r(k - \frac{n}{2})$ is the received signal and $x^*(k + \frac{n}{2})$ is the complex conjugate of the signal pattern we are looking for in the $n, f_d$ domain [12]. The procedure is done on blocks of signals, which have $N$ samples. Eq. 2 is a simplified version, where $y$ is the conjugate of $x$ multiplied by the rotational vector. If a distinct correlation peak is found in the function output, it means that the pattern matches the incoming signal with the given frequency shift and time delay. Over time, if the received signal contains the packets, the correlation peaks follow a tangent-function.

$$R(n, f_d) = \sum_{k=-N}^{N-1} r\left(k - \frac{n}{2}\right) \cdot x^*\left(k + \frac{n}{2}\right) e^{j2\pi f_d \frac{k}{N}} \quad (1)$$

$$R(n, f_d) = \sum_{k=-N}^{N-1} r\left(k - \frac{n}{2}\right) \cdot y\left(k + \frac{n}{2}, f_d\right) \quad (2)$$

Eq. 2 is a convolution, but one of the signals is reversed in time. Using the convolution theorem, we know that the product of the Fourier-transform of two signals is equal to the convolution of the two signals (Eq. 3). If we take the inverse Fourier-transform of the product (Eq. 4), we get the cross-correlation of the two signals [14].

$$\mathcal{F}\{R(n, f_d)\} = \mathcal{F}\left\{r\left(k - \frac{n}{2}\right)\right\} \cdot \mathcal{F}^*\left\{y\left(k + \frac{n}{2}, f_d\right)\right\} \quad (3)$$

$$R(\tau, f_d) = \mathcal{F}^{-1}\left[\mathcal{F}\left\{r\left(t - \frac{\tau}{2}\right)\right\} \cdot \mathcal{F}^*\left\{y\left(t + \frac{\tau}{2}, f_d\right)\right\}\right] \quad (4)$$

Since the signal processing is done on discrete time sampled signals, discrete Fourier-transform, more specifically Fast Fourier Transform (FFT) is used for better computation efficiency. To shift a signal in the frequency domain (to apply Doppler shift) the FFT of the signal is shifted circularly.

Fig. 2 shows the operations that are carried out periodically on the incoming signal. First, the reference signal is moved to frequency domain by FFT, which is shifted circularly so the frequency shifted copies of the signal without excessive computations are obtained. The incoming complex samples are moved to frequency domain by FFT, then the complex conjugate of the signal is multiplied by each frequency shifted

copy of the reference signal. After a two dimensional inverse FFT (IFFT) we get the double-sided delay-Doppler function.



Fig. 2. The signal processing chain

The output of the correlator is processed in blocks, whose size is determined by the length of the reference signal and the desired Doppler range. The complex float output is converted to absolute values on a logarithmic scale, then after a maximum search algorithm, the values are converted to a bitmap file. Lastly, the images are merged together into a montage, which yields the delay-Doppler image for the observed samples in blocks.

Fig. 3 shows the correlation peak at the center which represents the 0 delay and 0 Hz Doppler shift. As the incoming packets arrive at random times, the correlation peak moves left and right based on the sample delay and moves up and down according to the Doppler frequency. The length of the reference signal determines the correlation length and the maximum delay the function is able to show. On the vertical axis the range of Doppler shift is set manually, based on the maximum observable shift of the signal of the target satellite. In the case of Low Earth Orbit (LEO) satellite orbiting at 500 km above Earth the first orbital speed is defined as

$$v = \sqrt{\frac{Gm_E}{r}} =$$
$$= \sqrt{\frac{6.673 \cdot 10^{-11} \frac{Nm^2}{kg^2} \cdot 5.9 \cdot 10^{24} kg}{6370000m + 500000m}} = 7615 \frac{m}{s} \quad (5)$$

where $G$ is the the gravitational constant of Earth, $m_E$ is the mass of Earth and $r$ is the distance of the object from the center of the Earth. Using the velocity of the satellite and assuming a stationary receiver, we can calculate maximum observable Doppler frequency using Eq. 6, where $f_0$ is the transmitter frequency, $c$ is the speed of light and $v_t$ is the transmitter velocity.

$$f_D = \left| f_0 - f_0 \frac{c}{c + v_t} \right| = 11085 Hz \quad (6)$$

The resolution of the Doppler frequency estimation depends on the sampling rate ($f_s$) and the length of the FFT ($N$), the exact value is $\frac{f_s}{N}$.

## III. THE COMMUNICATION SUBSYSTEM OF MRC-100

Nowadays, small satellites use low cost, single chip radio transceivers to transfer data between the satellite and their ground station. Most of these use digital frequency modulation or frequency shift keying (FSK) for which a non-coherent receiver structure can be used. Therefore, bit and frame synchronization is done on a packet basis, using a preamble (usually 0101 sequences) and a sync word that has

Fig. 3. The delay-Doppler function with simulated signal, $0\,\mathrm{dB}$ SNR and $0\,\mathrm{Hz}$ Doppler

good auto-correlation properties. In the case of MRC-100 the communication subsystem (Fig. 4) uses an Acsip S68F integrated radio which is capable of handling LoRa and FSK transmissions (Fig. 6). It contains an SX1268 silicon along with a matching network that transforms the input/output impedance to $50\,\Omega$ and an integrated temperature compensated crystal oscillator (TCXO) as the reference clock.

As an essential component of the satellite, the communication module is designed to withstand single point failures by incorporating cold redundancy. Figure 5 shows how the two radios share the same antenna: we use radio frequency PIN (Positive Intrinsic Negative) diodes to connect them one at a time. In order to prevent over-current or overvoltage events, each radio is equipped with a dedicated limiter switch that automatically shuts off.

Telemetry packets use GMSK modulation and have a 32 bit standard 0101 preamble and a 32 bit sync word which is 0xE31C9DAE. This sequence is transmitted at the beginning of every single packet, so searching for this portion of data in the incoming signal regardless of the useful data in the packet produces cross-correlation peaks at the actual Doppler frequency.

## IV. MRC-100 TELEMETRY PROPERTIES

The mathematical formula of general frequency modulated (FM) signal is defined as

$$s_{FM}(t) = V_c cos\left(\omega_c t + 2\pi k_{FM}\int_0^t s_m(\tau)d\tau\right) \quad (7)$$

If we consider a sinusoidal modulating signal, $s_m(t) = V_m cos(\omega_m t)$ then the after integration Eq. 7 becomes

$$s_{FM}(t) = V_c cos\left(\omega_c t + 2\pi k_{FM}V_m\frac{sin(\omega_m t)}{\omega_m}\right) \quad (8)$$

Substituting $\omega_m$ by $2\pi f_m$ and defining $f_d = k_{FM}V_m$ as frequency deviation we get

$$s_{FM}(t) = V_c cos\left(\omega_c t + 2\pi f_d\frac{sin(2\pi f_m t)}{2\pi f_m t}\right) \quad (9)$$



Fig. 4. The UHF communication subsystem of MRC-100



Fig. 5. The block diagram of the UHF communication subsystem of MRC-100



Fig. 6. The block diagram of the Semtech SX1268 radio [7]

after simplification the the equation for an FM signal with sinusoidal modulating signal becomes

$$s_{FM}(t) = V_c cos \left( \omega_c t + \frac{f_d}{f_m} sin(2\pi f_m t) \right) \qquad (10)$$

If the modulating signal is a digital square wave, we get the resultant FSK signal (Eq. 11) by substituting the function of the sequence into Eq. 7, where $E_b$ is the bit energy, $T_b$ is the bit time, $\omega_c = 2\pi f_c$ is the carrier angular frequency, $k_{FM}$ is the modulation index, $d_k \in -1, +1$ is the data bit sequence and $P_{Tb}$ is the impulse function. If the data rate is

$$s_{FSK}(t) =$$
$$\sqrt{\frac{E_b}{T_b}} cos \left( \omega_c t + 2\pi k_{FM} \int_0^t \sum_{k=0}^{N-1} d_k P_{Tb}(\tau - kT_b) d\tau \right)$$
$$(11)$$

MRC-100 transmits telemetry packet groups (Fig. 7) periodically in normal operation. The period of the cycle is about 15 seconds and the default telemetry contains a LoRa (Long Range) modulated identification beacon, six GMSK telemetry packets and a synchronization packet. Each type has a different data rate: the synchronization packet is always transmitted at 5 kbps, the telemetry packets can be configured by ground station commands and their data rate is 12.5 kbps in normal operation. The LoRa beacon is also configurable, the default settings use SF12 with 12.5 kHz bandwidth with and air time of 892 ms [8].



Fig. 7. The default telemetry packets of MRC-100

### A. GMSK synchronization packet

MRC-100 uses a synchronization packet using GMSK modulation parameters and fixed data content that which is a pseudo random bit sequence (PRBS) sequence that helps the ground station schedule the telecommand packets.

By analyzing the the auto-correlation function of the packet, a peak that excels 12 dB from the noise floor can be seen (Fig. 8), meaning that a processing gain of 12 dB can be achieved. This makes detection possible even with a simple, stationary omnidirectional antenna with low gain, instead of using a directional antenna.

### V. SIMULATION RESULTS

The performance of the algorithm was evaluated using simulated signals. The time-domain samples of the frequency sync and frame sync bits are resampled, a square signal of +1 and -1 values are generated and filtered with a Gaussian window (BT=0.5) Finite Impulse Response (FIR) filter, which



Fig. 8. The normalized auto-correlation function of the RX sync packet

are then frequency modulated and multiplied by a carrier signal as shown in Fig. 9. The intermediate signals and the spectrum of the generated signal are shown in Fig. 10.



Fig. 9. Flowchart of the GMSK signal simulation



Fig. 10. Simulation of the GMSK modulated packet with highlighted bit sync and frame sync

The auto-correlation function of the simulated bit- and frame sync shown in Fig. 11 provides a 11 dB peak compared to the sidelobes, although the close vicinity of the peak is around 5 dB, due to the repeating 01 pattern.



Fig. 13. Doppler estimation for a simulated satellite pass with -15 dB SNR



Fig. 11. The normalized auto-correlation function of the bit and frame sync

For the input a list of Doppler frequencies were calculated using predict, which is an open-source software used for tracking satellites and celestial bodies [17]. We used the output to iteratively shift the generated signal in the frequency-domain so that it can be fed back to the Doppler estimator input.



Fig. 14. Error of Doppler estimation with -15 dB SNR

### A. Results using the bit- and frame sync samples as reference

Figures 12 and 13 show the delay-Doppler function for a satellite pass where the received signal to noise ratio (SNR) is 0 dB and -15 dB, respectively. Both results show distinct peaks that correspond with the Doppler shift of the input. The simulation outputs the measured frequency shift of the correlator output's maximum which was compared to the input values. Figure 14 shows the frequency error of the estimation. It should be noted that in the simulation the FFT bin size is $\frac{f_s}{n} = \frac{25000 Hz}{1024} = 24.41 Hz$ so the error range stayed within one FFT bin at every measurement point.



Fig. 12. Doppler estimation for a simulated satellite pass with 0 dB SNR

### B. Results using the synchronization packet

We also evaluated the algorithm using a collection of real-life samples recorded with a software defined radio (SDR) with zero Doppler shift. In this recording a variety of signals (Fig. 15) are present along with four RX synchronization packets ($2^{nd}, 4^{th}, 5^{th}$ and the last), to observe the behaviour of the algorithm. In this case the reference signal is longer, contains more energy so the correlation peak is narrower and more emphasized compared to the previous case where only the frame and bit synchronization was used as the reference signal. The correlator outputs the amplitude of the peak, which is proportional with the energy of the output signal according to Eq. **??**. This property of the correlator can be used to filter false detections. The amplitude threshold of the filter can be set by observing the distribution of estimated Doppler – correlation peak amplitude measurements. In Fig. 16 red and yellow markers show the valid and invalid measurement points which are classified by their correlation amplitude, shown in Fig. 17 where correlation peak value is displayed on a logarithmic scale. Yellow peaks typically have a larger amplitude than the red ones, indicating a possible detection of RX synchronization packet in the incoming signal. The point that shows about -300 Hz Doppler shift indicates a telecommand packet that was transmitted by the nearby ground station.

Picosatellite identification and Doppler estimation
using passive radar techniques



Fig. 15. Waterfall spectrum of several packets used for evaluation



Fig. 16. Delay-Doppler estimation of several packets used for evaluation



Fig. 17. The amplitude and Doppler shift of correlator peaks

The two results described above show that using the algorithm can be used to estimate the Doppler frequency of the incoming signal if the reference signal is present. However, the measurement output needs to be filtered, since the incoming signal does not always contain the pattern, so based on the amplitude of the correlator peaks, the invalid estimations can be discarded. Running the algorithm on a complete satellite pass can plot the Doppler curve of the orbit, which can be compared with the TLE used for tracking. This can be done with different TLEs that are candidates for the target satellite and the optimal TLE can be selected for future pass(es).

## VI. Conclusion

The work presented here described how Doppler estimation can be done using little information about the radio transmission parameters of a satellite that is too small to have a precise and up-to-date TLE, especially right after launch, when picosatellites are orbiting in a group with little distance to one another. Using cross-correlation and further signal processing, the Doppler frequency of the transmitted packets were estimated and invalid results were filtered, so an improved curve estimation can be done.

MRC-100 will be delivered to orbit with Transporter-8 Sun-Synchronous Orbit (SSO) Rideshare on 12 June, 2023. According to the launch broker, the two stage deployment will put the satellite 7-14 days after the rocket launch. After a successful beginning of the mission validation measurements will be carried out on actual signals received from the satellite.

## References

[1] Completed missions of Alba Orbital http://www.albaorbital.com/completed-missions Accessed October 25, 2023

[2] Unisat-7 launch campaign of GAUSS SRL https://www.gaussteam.com/satellites/gauss-latest-satellites/unisat-7/ Accessed October 25, 2023

[3] Firefly's first mission, "To The Black" https://fireflyspace.com/missions/flta002-to-the-black/ Accessed October 25, 2023

[4] Payloads of Transporter-3 mission https://www.eoportal.org/other-space-activities/transporter-3#payloads Accessed October 25, 2023

[5] Payloads of Transporter-5 mission https://www.eoportal.org/other-space-activities/transporter-5#passenger-payloads Accessed October 25, 2023

[6] Speretta, Stefano & Sundaramoorthy, Prem & Gill, Eberhard. (2017). "Long-term performance analysis of NORAD Two-Line Elements for CubeSats and PocketQubes.", *11th IAA Symposium on Small Satellites for Earth Observation,* Berlin, Germany

[7] Datasheet of Semtech SX1268 https://www.semtech.com/products/wireless-rf/lora-connect/sx1268# documentation Accessed October 25, 2023

[8] T. Herman and L. Dudás, "Satellite identification beacon system for PocketQube mission," *2022 24th International Microwave and Radar Conference (MIKON),* Gdansk, Poland, 2022, pp. 1–5, **DOI**: 10.23919/MIKON54314.2022.9924648.

[9] D. Mieczkowska et al., "Detection of objects on LEO using signals of opportunity," *2017 Signal Processing Symposium (SPSympo),* Jachranka, Poland, 2017, pp. 1–6, **DOI**: 10.1109/SPS.2017.8053660.

[10] M. Radmard, S. Bayat, A. Farina, S. Hajsadeghian and M. M. Nayebi, "Satellite-based forward scatter passive radar," *2016 17th International Radar Symposium (IRS),* Krakow, Poland, 2016, pp. 1–4, **DOI**: 10.1109/IRS.2016.7497275.

[11] V. I. Veremyev, E. N. Vorobev and Y. V. Kokorina, "Feasibility Study of Air Target Detection by Passive Radar Using Satellite-based Transmitters," *2019 IEEE Conference of Russian Young Researchers in Electrical and Electronic Engineering (EIConRus)*, Saint Petersburg and Moscow, Russia, 2019, pp. 154–157,
**DOI**: 10.1109/EIConRus.2019.8656630.

[12] T. Pető, L. Dudás and R. Seller, "DVB-T based passive radar," *2014 24th International Conference Radioelektronika*, Bratislava, Slovakia, 2014, pp. 1–4, https://doi.org/10.1109/Radioelek.2014.6828433.

[13] T. Pető, "Multichannel passive radar receiver platform," *2015 17th International Conference on Transparent Optical Networks (ICTON)*, Budapest, Hungary, 2015, pp. 1–4, **DOI**: 10.1109/ICTON.2015.7193445.

[14] Papoulis, A. "The Fourier Integral and Its Applications." New York: McGraw-Hill, pp. 244–245 and 252–253, 1962.

[15] Donát Takács, Boldizsár Markotics and Levente Dudás," Processing and Visualizing the Low Earth Orbit Radio Frequency Spectrum Measurement Results From the SMOG Satellite Project", *Infocommunications Journal*, Vol. XIII, No 1, March 2021, pp. 18–25.

[16] Yasir Ahmed Idris Humad, and Levente Dudás, "Wide Band Spectrum Monitoring System from 30MHz to 1800MHz with limited Size, Weight and Power Consumption by MRC-100 Satellite2, *Infocommunications Journal*, Vol. XIV, No 2, June 2022, pp. 56–63.,
**DOI**: 10.36244/ICJ.2022.2.6

[17] Homepage of predict. https://www.qsl.net/kd2bd/predict.html Accessed October 25, 2023.

**Tibor Herman** received his electrical engineering BSc and MSc degree from Budapest University of Technology and Economics in 2014 and 2016 respectively. He is currently a PhD student at the Department of Broadband Infocommunications an Electromagnetic Theory. His research area focuses on Small Satellite Subsystems, involving PocketQubes developed by students and lecturers of the department, radio communication hardware development and antenna design.

**Levente Dudás** received his electrical engineering MSc degree in 2007 and his PhD in 2018 at Budapest University of Technology and Economics. His research topic is Radar and Satellite Applications of Radio and Antenna Systems. He is currently a senior lecturer at the Department of Broadband Infocommunications an Electromagnetic Theory, working in the Radar and Remote Sensing Laboratory. His fields of interest are active and passive radar, CubeSat and PocketQube satellite development, analog high frequency hardware and antenna design, automated and remote controlled satellite tracking and signal processing.

# A Novel Strategic Caching and Availability Optimization for Wireless Unmanned Aerial Vehicle Communication Networks

Mohamed El Amrani, Hamid Garmani, Driss Ait Omar, Mohamed Ouaskou, Abdelkarim Ait Temghart,
and Mohamed Baslam

*Abstract*—The growing complexity in the 5G technology has created a necessity for UAV(Unmanned Aerial Vehicle)-assisted cellular networks as base stations. This is helpful for a wider coverage with higher transmission rates as it addresses three critical issues, i.e. location, performance and bandwidth. Besides, the integration of caching into conventional UAV infrastructure has received significant attentions since it can bring contents and memory storage closer to a mobile device. In a dynamic resource caching for wireless mobile networks, drone's settings contain important options for supporting a wide variety of applications and services, including the network access fee, quality of service (QoS) , number of cached contents, cache access fee and beaconing duration. A theoretic model based on game theory is developed to study the effect of competition among UAVs that have caching and sharing revenue model. Note that an optimal usage of UAV capabilities would thus lead to a cost-effective strategy for energy consumption and QoS requirements.

*Index Terms*—Wireless Communication, UAV, Beaconing Duration, Service fee, Quality of Service, Game Theory, Nash Equilibrium.

## I. INTRODUCTION

It is evident that drone applications are playing an increasingly prominent role in military, public, and civilian domains. Thus, this new technology needs to be further explored as they have been successfully used in various scenarios related to disaster and risk management, including earthquakes, landslides, floods, fires, tsunamis, etc ([1] [2]). Unmanned aerial vehicles (UAVs), better known as drones, are equipped with low-cost navigation sensors that enables detection, localization, and tracking of any target for more accurate predictions and the effectiveness of specific intervention options.

Not only can UAVs support effective both offensive and defensive military operations, but they can also be used as aerial base stations (BSs) to deliver reliable, cost-effective, and on-demand flying antenna systems. For the 5G network deployment in early stage, these artificial satellites relay and amplify radio telecommunication signals via a transponder so as to provide real-time communications with satellites enabling uplinks and downlinks ([3]). More precisely, UAV-assisted communications have several promising advantages compared to fixed nodes such as the flexibility of network deployment, real-time data, and high flexibility. Therefore,

UAV systems are expected to be the key component of the wireless because this platform can potentially facilitate mobile devices to be an all-coverage network to guarantee fast and ubiquitous connectivity in terrestrial cellular network.

Despite the progress made and the positive results achieved, drones are not able to launch and carry out a given mission with the autonomous flight capability. This is due to mainly technical limitations, including radio frequency (RF) restrictions, battery life, trajectory planning, severe weather conditions, restricted areas and moving obstacles. Besides, a number of logistical and privacy concerns are constantly rising when using drones as flying base stations. In order to gain a new perspective on this topic, the focus was placed on how to improve the efficiency of data transmission and optimize operational costs. In this context, we suggest UAV-aided small-cell content caching network to maximize the efficiency of content delivery and communicate locally near the end users. This process temporarily stores copies of files to use in high-speed memory chips, so that mobile devices can access the Internet content from drones to increase data retrieval performance. A cache's primary purpose is to extend coverage and provide high-quality network so as the requests for data can be served faster even if the origin server could be located anywhere in the world. In this context, Information-centric network offers a content centric paradigm to manage the explosion of content demand in the Internet architecture. ICN is a potential networking architecture for the Internet of Things. In ICN (Information-centric Networking), the content is requested by using unique names instead of IP, every node in the network can cache and serve the requested content. Benefits of ICN are a robust communication for ad-hoc networks, better response time using network caching, improved bandwidth utilization, security, mobility, etc ([4]).

In this paper, we use a new economic game theory model in which the autonomous coordinated flying for groups of UAVs can be deployed in real time to maximize network coverage. In this case, a non-cooperative game between drones is formulated to meet some target quality of service (QoS) while having constraints specified by beaconing duration and service fee. The proposed model of the competitive game among drones attempt to reach a Nash equilibrium point. In this steady state, there is no incentive for all drones to deviate from their strategies for changing their outcomes. Moreover, we design a distributed algorithm that converges to a situation representing the optimal solution. This is a concern that may

be best addressed using Bertrand game model to derive the fee of fairness and then analyze the equilibrium distribution of outputs among the drones. Numerical experiments are conducted to seek and evaluate the factors that impact the drones' strategies on their expected profits.

The remainder of this paper is organized as follows. In Section 2, we first discuss briefly some related works. The theoretical models presented in Section 3 aim specifically to formulate and model UAV systems and their utility function that explain drones' behavior among risky or uncertain choices. In the same line, Section 4 and Section 5 provide theoretical analysis of the considered non-cooperative game theory that focuses on the Nash equilibrium of the mixed-strategy best-response. The results of the simulation are reported through Section 6, with concluding discussion in Section 7.

## II. RELATED WORKS

There have been several recent studies where UAV-assisted cellular networks is proposed as a cost-effective solution for ubiquitous coverage. The authors in [5] [6] study cellular networks that implement cellular radio access nodes to meet the requirements of 5G New Radio (NR) performance. In this context, game theory is used to address more complicated problems of optimizing resource allocation in UAV-based communication. In [7] the authors introduced a comprehensive review of the existing game theoretic techniques that handle various applications of drone-based communication networks. The authors in [8] used a non-cooperative sub-modular game to model beaconing periods scheduling manner, to maximize the coverage probability of mobile devices. To support ground-based units, the authors in [9] rely on Genetic algorithms and non-cooperative games for ensuring the optimal flying solutions and maximizing coverage as well. In [10], the authors proposed a theoretical framework to model the interaction of UAVs that act as a flying base station. To this aim, they use a non-cooperative game model to determine the best pricing strategy while guaranteeing high levels of UAVs availability. Due to the limited on-board battery size, it is necessary to define optimal periodic beaconing by taking into account the limited battery capacity of the UAVs and their difficulties in recharging. Besides, edge-caching has received much attention as an efficient technique to reduce the latency to access popular content and overcome backhaul congestion, especially during peak traffic [11] [12]. Similarly, the authors in [13] presented a novel proposed scheme based on proactive caching to support the drone with limited flight endurance and payload capacity. More precisely, the proposed solution aims to minimize the file caching cost and recovery cost. This was achieved mainly by jointly optimizing the drone communication scheduling, drone trajectory, and file caching policy. In addition to focusing on the overall network performance, the authors in [14] have suggested a new scheme to ensure secure transmission for UAV-relayed wireless networks with caching capability. In [15] the authors implemented an online caching-based wireless UAV by jointly optimizing UAV trajectory, transmission power, and caching scheduling. For an effective architecture that improves the wireless coverage, authors in [16] employ backscattering

communication (BackCom) to transmit data to guarantee the effectiveness of real-time transmission and minimize the data collection latency. This will undoubtedly contribute to the further strengthening the UAV's lifetime network and improve the ground cellular networks' coverage as well.

As outlined above, the activity scheduling of the UAV as aerial base stations has been extensively investigated during the past few years. Those studies, however, have mainly examined various strategic decisions to define the optimal fee and the appropriate beaconing duration of UAVs with limited battery capacity. Thus, this work is specifically carried out to explore fair competition between UAVs having caching and sharing revenue model associated with energy efficiency improvements. Furthermore, to the best of our knowledge, none of the previous work analyses the most significant impacts of using caching service and sharing revenue model in a UAVs network. To explore these issues, the present paper moves toward this less explored case, where each UAV chooses the network access fee, QoS, number of cached contents, cache access fee and beaconing duration. According to this specific situation, we conduct simulation experiments to demonstrate the practicality of our approach and show how caching and sharing revenue model affect the UAV's energy efficiency as well as the QoS and the pricing strategies.

## III. PROBLEM FORMULATION

In this paper, we consider a telecommunication network having G UAVs, in which each UAV is in competition with the other UAVs for the users on the ground. The monetary flow between different entities is shown in Figure 1 with different fees as described in the Table I. Each UAV $j$ picks its availability duration $\xi_j$ represented by the periodic beaconing time chosen within the interval $[0, T]$, a service fee per data unit $p_{s_j}$, a content access fee $p_{c_j}$, QoS $q_{s_j}$ and number of cached items $H_j = \sum_{f=1}^{F} h_{jf}$.

### A. Notations



Fig. 1. Model architecture.

### B. Service probability

The UAVs are connected through wireless backhaul to the core network and move randomly according to a Random

A Novel Strategic Caching and Availability Optimization for
Wireless Unmanned Aerial Vehicle Communication Networks

TABLE I
SUMMARY OF NOTATION.

| Notation | Description |
|---|---|
| $G$ | Number of UAVs. |
| $F$ | Number of items. |
| $p_{s_j}$ | Network access fee of $UAV_j$. |
| $p_{c_j}$ | Content access fee of $UAV_j$. |
| $q_{s_j}$ | Quality of service of $UAV_j$. |
| $\xi_j$ | Beaconing period durations of $UAV_j$. |
| $\mu_j$ | Caching cost of $UAV_j$. |
| $\vartheta_j$ | backhaul bandwidth cost. |
| $T$ | Time slot. |
| $m$ | Time window. |
| $p_{t_j}$ | Transmission fee paid by $UAV_j$. |
| $\rho_j^g$ | Sensitivity of $UAV_j$ to fee $p_{s_g}$ of $UAV_g$. |
| $\sigma_j^g$ | Sensitivity of $UAV_j$ to QoS $q_{s_g}$ of $UAV_g$. |
| $\alpha_j^g$ | Sensitivity of $UAV_j$ to fee $p_{c_g}$ of $UAV_g$. |
| $\delta_j^g$ | Sensitivity of $UAV_j$ to beaconing period $\xi_g$ of $UAV_g$. |
| $P_{srv_j}$ | Successful contact probability of $UAV_j$. |
| $D_j^0$ | the potential demand of ground users of $UAV_j$. |
| $d_j$ | Demand of $UAV_j$. |
| $B_j$ | Backhaul bandwidth. |
| $f^\eta$ | the rank of item $f$. |
| $\eta$ | the skewness of the popularity distribution. |
| $C_{b_j}$ | The energy cost for sending beacons of $UAV_j$. |
| $C_{q_j}$ | is the energy cost for providing QoS of $UAV_j$. |
| $C_{h_j}$ | is the energy cost for caching content of $UAV_j$. |
| $C_{s_j}$ | is the energy consumed to switch the state of transceiver of $UAV_j$. |
| $CP$ | Content Provider. |
| $UAV$ | Unmanned Aerial Vehicle. |

Way-point mobility model to cover a specific area. Each UAVs send a beacon to the ground users to announce its presence during a specific period of duration $\xi$ . The UAVs choose their beaconing period durations to maximize the probability of encountering a mobile device on the ground and send periodic beacon advertising his availability for users on the ground. The beacon/idle cycle is periodically repeated every time slot $T$ during a time window: $m = L \times T$. However UAVs should define their beaconing periods strategically in order to maximize their encounter rate with the ground users. They should avoid battery depletion resulting from maintaining useless beaconing in the absence of contact opportunities. The first encounter follows an exponential distribution with a random parameter $\lambda$. In order for a UAV $j$ to encounter first the ground users at time $\xi$, the following conditions must hold: the UAV has to be beaconing at $\xi$. The encounters need to be happen while UAV $j$ competitors are inactive. In other words, all encounters of other UAVs happen before the times instant $\xi_j$ must be unsuccessful. Consequently, the successful contact probability is given by the following equation:

$$P_{srv_j} = \sum_{g=1}^{G} \left[ P(T_j \leq T_g) + P(T_j \geq T_g) P_{slp_g} \right] P_{bcn_j} \quad (1)$$

We define the probability of $UAV_j$ beaconing while encountering for the first time the destination within $[0, m]$ :

$$P_{bcn_j} = \sum_{s=0}^{l-1} \left( \int_{sT}^{sT+\xi_j} \lambda_j e^{-\lambda_j x} dx \right) \quad (2)$$

$$= -\frac{e^{-\lambda T}(e^{-m\lambda_j} - e^{-\lambda_j(m+\xi_j)} - 1 + e^{-\lambda_j \xi_j})}{e^{\lambda_j T} - 1}$$

For a UAV $j$, the probability of being idle is given by the following equation:

$$P_{slp_j} = \sum_{s=0}^{l-1} \left( \int_{sT+\xi_j}^{(s+1)T} \lambda_j e^{-\lambda_j x} dx \right) \quad (3)$$

$$= \frac{e^{-\lambda_j T}(-e^{-\lambda_j(m+\xi_j)} + e^{-\lambda_j(m+T)} + e^{-\lambda_j \xi_j} - e^{-\lambda_j T})}{e^{\lambda_j T} - 1}$$

The probability that UAV $j$ encounters first the ground destination without accounting for its state (probing/idle) is expressed as follows:

$$P(T_j \leq T_g) = \frac{\lambda_g e^{-m(\lambda_g + \lambda_j)} + (-\lambda_g - \lambda_j)e^{-\lambda_g m} + \lambda_j}{\lambda_j + \lambda_g} \quad (4)$$

And finally, we define the probability that UAV $j$ encounters first the ground destination without accounting for its state:

$$P(T_j \geq T_g) = \frac{\lambda_j e^{-m(\lambda_g + \lambda_j)} + (-\lambda_g - \lambda_j)e^{-\lambda_j m} + \lambda_g}{\lambda_j + \lambda_g} \quad (5)$$

### C. Demand model

The ground users demand are affected by three market parameters service fee, beaconing duration and QoS. The linear demand function is [17][18][19]:

$$\begin{aligned} d_j = &D_j^0 - \rho_j^j p_{s_j} + \delta_j^j \xi_j + \alpha_j^j p_{c_j} + \sigma_j^j q_{s_j} \\ &+ \sum_{g=1, g \neq j}^{G} (\rho_j^g p_{s_g} - \delta_j^g \xi_g + \alpha_j^g p_{c_g} - \sigma_j^g q_{s_g}) \end{aligned} \quad (6)$$

The parameter $D_j^0$ expresses the potential demand of ground users. $\rho_j^g$, $\delta_j^g$, $\alpha_j^g$ and $\sigma_j^g$ are positive parameters representing respectively the responsiveness of $UAV_j$ to fee $p_{s_g}$, beaconing period $\xi_g$, fee $p_{c_g}$ and QoS $q_{s_g}$ of $UAV_g$. For $UAV_j$, the demand $d_j$ is decreasing in the fee it charges, $p_{s_j}$, $p_{c_g}$, and increase in the fee charged by its opponent, $p_{s_g}$, $p_{c_g}$, $g \neq j$. The analogous relationship holds in QoS and beaconing period, in this case $D_j$ is increasing in $q_{s_j}$ (resp. $\xi_j$) and decreasing in $q_{s_g}$ (resp. $\xi_g$).

**Assumption 1** *The sensitivity $\rho$ verifies:*

$$\rho_j^j \geq \sum_{g=1, g \neq j}^{G} \rho_j^g$$

*The sensitivity $\delta$ verifies:*

$$\delta_j^j \geq \sum_{g=1, g \neq j}^{G} \delta_j^g$$

*The sensitivity $\sigma$ verifies:*

$$\sigma_j^j \geq \sum_{g=1, g \neq j}^{G} \sigma_j^g$$

*The sensitivity $\alpha$ verifies:*

$$\alpha_j^j \geq \sum_{g=1, g \neq j}^{G} \alpha_j^g$$

Assumption 1 will be needed to ensure the uniqueness of the resulting equilibrium. It is furthermore a reasonable condition, in that Assumption 1 implies that the influence of an UAV fee (resp. beaconing duration) is significantly greater on its observed demand than the fees of its opponents.

### D. Utility function

The utility $U_j$ of $UAV_j$ is the difference between the obtained reward and the associated costs

$$
\begin{aligned}
U_j = \sum_{f=1}^{F} \Theta_f \{ &(p_{s_j} - p_{t_j})(1 - h_{jf})P_{srv_j}d_j \\
&+ (p_{s_j} + p_{c_j} - \mu_j)h_{jf}P_{srv_j}d_j\} \\
&- \frac{(C_{b_j}\xi_j + C_{q_j}q_{s_j} + C_{h_j}\sum_{f=1}^{F}h_{jf} + C_{s_j})l}{m} \\
&- \vartheta_j(F - \sum_{f=1}^{F}h_{jf})B_j
\end{aligned}
\tag{7}
$$

$\sum_{f=1}^{F}\Theta_f(p_{s_j} - p_{t_j})(1 - h_{jf})P_{srv_j}d_j$ is the revenue of $UAV_j$ by serving the request demand $\sum_{f=1}^{F}\Theta_f(1 - h_{jf})P_{srv_j}d_j$. $\mu_j \sum_{f=1}^{F}\Theta_f h_{jf}P_{srv_j}d_j$ is the caching fee paid by the $UAV_j$ when serving the demand $\sum_{f=1}^{F}\Theta_f h_{jf}P_{srv_j}d_j$ of the item $f$ from its cache. Each ground user requests an item, which is selected independently according to a discrete distribution $\Theta_f$ where $1 \leq f \leq F$, and $F$ is the library size. We assume that the item $f$ is requested by their popularity, characterized by Zipf popularity distribution $\Theta_f$ [20] [21] [22]. The Zipf popularity distribution of item $f$ is defined by $\Theta_f = A^{-1}f^{-\eta}$, where $A = \sum_{f=1}^{F}f^{-\eta}$, $f^\eta$ is the rank of item $f$, and $\eta$ is the skewness of the popularity distribution. $\frac{(C_{b_j}\xi_j + C_{q_j}q_{s_j} + C_{h_j}\sum_{f=1}^{F}h_{jf} + C_{s_j})l}{m}$ is the energy consumed. $\vartheta_j(F - \sum_{f=1}^{F}h_{jf})B_j$ is a fee paid by $UAV_j$. $B_j$ is the backhaul bandwidth required by the $UAV_j$. The backhaul bandwidth $B_j$ of $UAV_j$ is expressed as [23] [24] [25]

$$B_j = (F - \sum_{f=1}^{F}h_{jf})(P_{srv_j}d_j + q_{s_j}^2) \tag{8}$$

Then, the utility function is given by

$$
\begin{aligned}
U_j = \sum_{f=1}^{F} \Theta_f \{ &(p_{s_j} - p_{t_j})(1 - h_{jf})P_{srv_j}d_j \\
&+ (p_{s_j} + p_{c_j} - \mu_j)h_{jf}P_{srv_j}d_j\} \\
&- \frac{(C_{b_j}\xi_j + C_{q_j}q_{s_j} + C_{h_j}\sum_{f=1}^{F}h_{jf} + C_{s_j})l}{m} \\
&- \vartheta_j(F - \sum_{f=1}^{F}h_{jf})(P_{srv_j}d_j + q_{s_j}^2)
\end{aligned}
\tag{9}
$$

sectionGame analysis Let $\psi = [\mathscr{G}, \{P_{s_j}, Q_{s_j}, \Xi_j, P_{c_j}\}, \{U_j(.)\}]$ denote the non-cooperative game (NPQBPG), where $\mathscr{G}=\{1, ..., G\}$ is the set of UAVs, $P_{s_j}$ is the network access fee strategy set of $UAV_j$, $Q_{s_j}$ is the QoS strategy set of $UAV_j$, $\Xi_f$ is the beconing strategy set of $UAV_j$, $P_{c_j}$ is the content access fee strategy set of $UAV_j$, and $U_j(.)$ is the utility function of $UAV_j$. We assume that the strategy spaces $P_{s_j}, Q_{s_j}, \Xi_j$ and $P_{c_j}$ of each $UAV_j$ are compact and convex sets with maximum and minimum constraints. Thus, for each $UAV_j$ we consider as respective strategy spaces the closed intervals: $P_{s_j} = \left[\underline{p}_{s_j}, \overline{p}_{s_j}\right]$, $Q_{s_j} = \left[\underline{q}_{s_j}, \overline{q}_{s_j}\right]$, $\Xi_j = \left[\underline{\xi}_j, \overline{\xi}_j\right]$ and $P_{c_j} = \left[\underline{p}_{c_j}, \overline{p}_{c_j}\right]$. Let the fee vector $p_s = (p_{s_1}, ..., p_{s_G})^T \in P_s^G = P_{s_1} \times P_{s_2} \times ... \times P_{s_G}$, QoS vector $q_s = (q_{s_1}, ..., q_{s_G})^T \in Q_s^G = Q_{s_1} \times Q_{s_2} \times ... \times Q_{s_G}$, beconing vector $\Xi = (\xi_1, ..., \xi_G)^T \in \Xi^G = \Xi_1 \times \Xi_2 \times ... \times \Xi_G$, fee vector $p_c = (p_{c_1}, ..., p_{c_G})^T \in P_c^G = P_{c_1} \times P_{c_2} \times ... \times P_{c_G}$,

### E. Fee game

A $NPQBCG$ in fee $p_s$ is defined for fixed $\mathbf{q}_s \in Q_s, \xi \in \Xi$, $\mathbf{p}_c \in P_c$ as $\Psi(\mathbf{q}_s, \xi, \mathbf{p}_c) = [G, \{P_{s_j}\}, \{U_j(., \mathbf{q}_s, \xi, \mathbf{p}_c)\}]$.

**Definition 1** *A fee vector $\mathbf{p_s}^* = (p_{s_1}^*, ..., p_{s_G}^*)$ is a Nash equilibrium of the NPQBG $\Psi(\mathbf{q}_s, \xi, \mathbf{p}_c))$ if:*

$$\forall (j, p_{s_j}) \in (G, P_{s_j}),$$
$$U_j(p_{s_j}^*, \mathbf{p}_{s_{-j}}^*, \mathbf{q}_s, \xi, \mathbf{p}_c) \geq U_j(p_{s_j}, \mathbf{p}_{s_{-j}}^*, \mathbf{q}_s, \xi, \mathbf{p}_c) \tag{10}$$

**Theorem 1** *For each $\mathbf{q}_s \in Q_s$, $\xi \in \Xi$, $\mathbf{p}_c \in P_c$ the game $[G, \{P_{s_j}\}, \{U_j(., \mathbf{q}_s, \xi, \mathbf{p}_c)\}]$ admit a unique Nash equilibrium.*

$$\frac{\partial^2 U_j}{\partial p_{s_j}^2} = -2\rho_j^j P_{srv_j} \leq 0 \tag{11}$$

The second derivative of the utility function is negative, then the utility function is thus concave, which ensures existence of a Nash equilibrium point in the game $\Psi(\mathbf{q}_s, \xi, \mathbf{p}_c)$.

We use the following proposition that holds for a concave game [26]: If a concave game satisfies the dominance solvability condition :

$$-\frac{\partial^2 U_j}{\partial p_{s_j}^2} \geq \sum_{g=1, g \neq j}^{G} \left| \frac{\partial^2 U_j}{\partial p_{s_j} \partial p_{s_g}} \right|. \tag{12}$$

then the game $\Psi(\mathbf{q}_s, \xi, \mathbf{p}_c)$ admits a unique Nash equilibrium.
The mixed partial derivative is written as:

$$\frac{\partial^2 U_j}{\partial p_{s_j} \partial p_{s_g}} = \rho_j^g P_{srv_j} \tag{13}$$

Then,

$$-\frac{\partial^2 U_j}{\partial p_{s_j}^2} - \sum_{g=1,g\neq j}^{G} \left|\frac{\partial^2 U_j}{\partial p_{s_j}\partial p_{s_g}}\right| = P_{srv_j}(\rho_j^j - \sum_{g=1,g\neq j}^{G} \rho_j^g) \geq 0$$

(14)

Thus, the game $\Psi(\mathbf{q}_s, \xi, \mathbf{p}_c)$ admits a unique Nash equilibrium point.

*F. Beaconing duration game*

A $NPQBCG$ in beaconing duration $\xi$ is defined for fixed $\mathbf{p}_s \in P_s$, $\mathbf{q}_s \in Q_s$ and $\mathbf{p}_c \in P_c$ as $\Psi(\mathbf{p}_s, \mathbf{q}_s, \mathbf{p}_c) = [G, \{\xi_j\}, \{U_j(\mathbf{p}_s, \mathbf{q}_s, ., \mathbf{p}_c)\}]$.

**Definition 2** *A beaconing vector $\xi^* = (\xi_1^*, ..., \xi_G^*)$ is a Nash equilibrium of the $NPQBCG$ $\Psi(\mathbf{p}_s, \mathbf{q}_s, \mathbf{p}_c)$ if:*

$$\forall(j, \xi_j) \in (G, \Xi_j),$$
$$U_j(\mathbf{p}_s, \mathbf{q}_s, \xi_j^*, \xi_{-j}^*, \mathbf{p}_c) \geq U_j(\mathbf{p}_s, \mathbf{q}_s, \xi_j, \xi_{-j}^*, \mathbf{p}_c)$$

(15)

**Theorem 2** *For each $\mathbf{q}_s \in Q_s$, $\mathbf{p}_s \in P_s$ and $\mathbf{p}_c \in P_c$, the game $[G, \{\xi_j\}, \{U_j(\mathbf{p}, \mathbf{q}_s, .)\}]$ admit a unique Nash equilibrium.*

$$\frac{\partial^2 U_j}{\partial \xi_j^2} = (u_j - \nu_j(F - \sum_{f=1}^{F} h_{jf}))$$
$$\times \frac{e^{-\lambda_j(\xi_j+T)}(e^{-\lambda_j m} - 1)(-2\delta_j^i + d_j\lambda_j)}{e^{\lambda_j T} - 1}$$

(16)

where $u_j = \sum_{f=1}^{F} \Theta_f\{(p_{s_j} - p_{t_j})(1 - h_{jf}) + (p_{s_j} + p_{c_j} - \mu_j)h_{jf}\}$.

we assume $u_j \geq \nu_j(F - \sum_{f=1}^{F} h_{jf})$ and $e^{-\lambda_j m} \leq 1$ then,

$$\frac{\partial^2 U_j}{\partial \xi_j^2} \leq 0$$

(17)

The second derivative of the utility function is negative, then the utility function is thus concave, which ensures existence of a Nash equilibrium in the game $\Psi(\mathbf{p}_s, \mathbf{q}_s, \mathbf{p}_c)$.

We use the following proposition that holds for a concave game [26]: If a concave game satisfies the dominance solvability condition :

$$-\frac{\partial^2 U_j}{\partial \xi_j^2} \geq \sum_{g=1,g\neq j}^{G} \left|\frac{\partial^2 U_j}{\partial \xi_j\partial \xi_g}\right|$$

(18)

then the game $\Psi(\mathbf{p}_s, \mathbf{q}_s, \mathbf{p}_c)$ admits a unique Nash equilibrium point.

The mixed partial is written as:

$$\frac{\partial^2 U_j}{\partial \xi_j\partial \xi_g} = -\frac{P(T_j > T_n)\lambda_j e^{-\lambda_j T}(e^{\lambda_j m} - 1)\lambda_g e^{\lambda_g T}}{(e^{\lambda_j T} - 1)}$$
$$\times (e^{-\lambda_g m} - 1)(u_j - \nu_j(F - \sum_{f=1}^{F} h_{jf}))d_j e^{-\lambda_j \xi_j - \lambda_g \xi_g}$$
$$+ \frac{\lambda_j e^{-\lambda_j T}(e^{\lambda_j m} - 1)\delta_j^g(u_j - \nu_j(F - \sum_{f=1}^{F} h_{jf}))e^{-\lambda_j \xi_j}}{e^{\lambda_j T} - 1}$$
$$+ \frac{\delta_j^j(u_j - \nu_j(F - \sum_{f=1}^{F} h_{jf}))\lambda_g e^{\lambda_g T}(e^{-\lambda_g m} - 1)e^{\lambda_g \xi_n}}{e^{-\lambda_g} - 1}$$

(19)

then,

$$-\frac{\partial^2 U_j}{\partial \xi_j^2} - \sum_{g=1,g\neq j}^{G} \left|\frac{\partial^2 U_j}{\partial \xi_j\partial \xi_g}\right| = \frac{\lambda_j e^{-\lambda_j(T+\xi_j)}(e^{-\lambda_j m} - 1)}{e^{\lambda_j T} - 1}$$
$$\times (u_j - \nu_j(F - \sum_{f=1}^{F} h_{jf}))(2\delta_j^j - \sum_{g=1,g\neq j}^{G} \delta_j^g)$$
$$+ \frac{e^{-\lambda_j \xi_j}(e^{-\lambda_j m} - 1)}{e^{\lambda_j T} - 1}u_j - \nu_j(F - \sum_{f=1}^{F} h_{jf})d_j$$
$$\times (-\lambda_j^2 e^{-\lambda_j \xi_j} + \lambda_j e^{-\lambda_j T}\sum_{g=1,g\neq j}^{G} P(T_j > T_g)\lambda_n$$
$$\times (e^{-\lambda_g m} - 1)e^{-\lambda_g \xi_g}\frac{e^{\lambda_g T}}{e^{\lambda_g} - 1}$$
$$+ \delta_j^j(u_j - \nu_j(F - \sum_{f=1}^{F} h_{jf}))(\frac{\lambda_j e^{-\lambda_j(T+\xi_j)}(e^{-\lambda_j m} - 1)}{e^{\lambda_j T} - 1}$$
$$- \sum_{g=1,g\neq j}^{G} \frac{\lambda_g e^{\lambda_g(T-\xi_g)}(e^{-\lambda_g m} - 1)}{e^{\lambda_g} - 1}) \geq 0$$

(20)

Thus, the game $\Psi(\mathbf{p}_s, \mathbf{q}_s, \mathbf{p}_c)$ admits a unique Nash equilibrium point.

*G. Quality of service game*

A $NPQBCG$ in QoS is defined for fixed $\mathbf{p}_s \in P_s$, $\xi \in \Xi$ and $\mathbf{p}_c \in P_c$ as $\Psi(\mathbf{p}_s, \xi, \mathbf{p}_c) = [G, \{Q_{s_j}\}, \{U_j(\mathbf{p}_s, ., \xi, \mathbf{p}_c)\}]$.

**Definition 3** *A QoS vector $\mathbf{q}_s^* = (q_{s_1}^*, ..., q_{s_G}^*)$ is a Nash equilibrium of the $NPQBCG$ $\Psi(\mathbf{p}_s, \xi, \mathbf{p}_c)$ if*

$$\forall(j, q_{s_j}) \in (G, Q_{s_j}),$$
$$U_j(\mathbf{p}_s, q_{s_j}^*, \mathbf{q}_{s-j}^*, \xi, \mathbf{p}_c) \geq U_j(\mathbf{p}_s, q_{s_j}, \mathbf{q}_{s-j}^*, \xi, \mathbf{p}_c)$$

**Theorem 3** *For each $\mathbf{p} \in P_s$, $\xi \in \Xi$ and $\mathbf{p}_c \in P_c$ the game $[G, \{Q_{s_j}\}, \{U_j(\mathbf{p}, ., \xi, \mathbf{p}_c)\}]$ admits a unique Nash equilibrium.*

$$\frac{\partial^2 U_j}{\partial q_{s_j}^2} = -2\nu_j(F - \sum_{f=1}^{F} h_{jf}) \leq 0$$

(21)

the utility function is concave, which ensures existence of a Nash equilibrium point in the game $\Psi(\mathbf{p}, \xi, \mathbf{p}_c)$.

In order to prove uniqueness, we follow, [27], and define the weighted sum of user utility functions.

$$\psi(\mathbf{q}_s, \mathbf{x}) = \sum_{j=1}^{G} x_j U_j(q_{s_j}, \mathbf{q}_{s-j})$$

(22)

The pseudo-gradient of (22) is given by:

$$v(\mathbf{q}_s, \mathbf{x}) = \left[x_1 \nabla U_1(q_{s_1}, q_{s-1}), ..., x_G \nabla U_G(q_{s_G}, \mathbf{q}_{s-G})\right]^T$$

(23)

The Jacobian matrix J of the pseudo-gradient (w.r.t.q) is written:

$$
J = \begin{pmatrix}
x_1 \frac{\partial^2 U_1}{\partial q_{s_1}^2} & x_1 \frac{\partial^2 U_1}{\partial q_{s_1}\partial q_{s_2}} & \cdot & \cdot & \cdot & x_1 \frac{\partial^2 U_1}{\partial q_{s_1}\partial q_{s_G}} \\
x_2 \frac{\partial^2 U_2}{\partial q_{s_2}\partial q_{s_1}} & x_2 \frac{\partial^2 U_2}{\partial q_{s_2}^2} & \cdot & \cdot & \cdot & x_2 \frac{\partial^2 U_2}{\partial q_{s_2}\partial q_{s_G}} \\
\cdot & \cdot & \cdot & & & \cdot \\
\cdot & \cdot & & \cdot & & \cdot \\
\cdot & \cdot & & & \cdot & \cdot \\
x_G \frac{\partial^2 U_G}{\partial q_{s_G}\partial q_{s_1}} & x_G \frac{\partial^2 U_G}{\partial q_{s_G}\partial q_{s_2}} & \cdot & \cdot & \cdot & x_G \frac{\partial^2 U_G}{\partial q_{s_G}^2}
\end{pmatrix}
$$

$$
= \begin{pmatrix}
\Lambda_1 & 0 & \cdot & \cdot & \cdot & 0 \\
0 & \Lambda_2 & \cdot & \cdot & \cdot & 0 \\
\cdot & \cdot & \cdot & & & \cdot \\
\cdot & \cdot & & \cdot & & \cdot \\
\cdot & \cdot & & & \cdot & \cdot \\
0 & 0 & \cdot & \cdot & \cdot & \Lambda_G
\end{pmatrix}
$$

where $\Lambda_i = -2x_i\nu_i(F - \sum_{f=1}^{F} h_{if})$, $i = 1, ..., G$.

Thus, $J$ is a diagonal matrix with negative diagonal elements. This implies that $J$ is negative definite. Henceforth $[J + J^T]$ is also negative definite, and according to Theorem (6) in, [27], the weighted sum of the utility functions $\psi(q_s, x)$ is diagonally strictly concave. Thus, the game $G(\mathbf{p}, \xi)$ admits a unique Nash equilibrium point is unique.

### H. Fee $P_c$ game

A $NPQBCG$ in fee $p_c$ is defined for fixed $\mathbf{p}_s \in P_s$, $\mathbf{q}_s \in Q_s$ and $\xi \in \Xi$ as $\Psi(\mathbf{p}_c, \mathbf{q}_s, \xi) = [G, \{P_{s_j}\}, \{U_j(\mathbf{p}_s, \mathbf{q}_s, \xi, .)\}]$.

**Definition 4** *A fee vector* $\mathbf{p_c}^* = (p_{c_1}^*, ..., p_{c_G}^*)$ *is a Nash equilibrium of the* $NPQBCG$ $\Psi(\mathbf{p}_s, \mathbf{q}_s, \xi)$ *if:*

$$
\forall (j, p_{c_j}) \in (G, P_{c_j}),
$$
$$
U_j(\mathbf{p}_s, \mathbf{q}_s, \xi, p_{c_j}^*, \mathbf{p}_{c-j}^*) \geq U_j(\mathbf{p}_s, \mathbf{q}_s, \xi, p_{c_j}, \mathbf{p}_{c-j}^*) \quad (24)
$$

**Theorem 4** *For each* $\mathbf{p}_s \in P_s$, $\mathbf{q}_s \in Q_s$ *and* $\xi \in \Xi$ *the game* $[G, \{P_{s_j}\}, \{U_j(\mathbf{p}_s, \mathbf{q}_s, \xi, .)\}]$ *admit a unique Nash equilibrium.*

$$
\frac{\partial^2 U_j}{\partial p_{c_j}^2} = -2\alpha_j^j P_{srv_j} \sum_{f=1}^{F} \Theta_f h_{jf} \leq 0 \quad (25)
$$

The second derivative of the utility function is negative, then the utility function is thus concave, which ensures existence of a Nash equilibrium point in the game $\Psi(\mathbf{q}_s, \xi, \mathbf{p}_c)$.

We use the following proposition that holds for a concave game [26] : If a concave game satisfies the dominance solvability condition :

$$
-\frac{\partial^2 U_j}{\partial p_{c_j}^2} \geq \sum_{g=1, g \neq j}^{G} \left| \frac{\partial^2 U_j}{\partial p_{c_j} \partial p_{c_g}} \right| \quad (26)
$$

then the game $\Psi(\mathbf{p}_s, \mathbf{q}_s, \xi)$ admits a unique Nash equilibrium.

The mixed partial is written as:

$$
\frac{\partial^2 U_j}{\partial p_{c_j} \partial p_{c_g}} = \alpha_j^g P_{srv_j} \sum_{f=1}^{F} \Theta_f h_{jf} \quad (27)
$$

Then,

$$
-\frac{\partial^2 U_j}{\partial p_{c_j}^2} - \sum_{g=1, g \neq j}^{G} \left| \frac{\partial^2 U_j}{\partial p_{c_j} \partial p_{c_g}} \right| = \left(2\alpha_j^j - \sum_{g=1, g \neq j}^{G} \alpha_j^g\right)
$$
$$
\times P_{srv_j} \sum_{f=1}^{F} \Theta_f h_{jf} \geq 0 \quad (28)
$$

Thus, the game $\Psi(\mathbf{p}_s, \mathbf{q}_s, \xi)$ admits a unique Nash equilibrium point.

### I. Learning Nash equilibrium

The problem presented in this document respects the uniqueness of equilibrium of NASH. After step Nash's uniqueness of equilibrium comes the second step on how to design an algorithm that converges to this equilibrium.

The fundamentals of the best dynamic response schemes, which can lead to a Nash equilibrium can be represented according to the following description: Let $G$ be a non-cooperative strategic game. Maximizing utility by UAV's response strategy by considering the strategies of other UAVs is the best. The importance of best response is useful if the game converges to a stable state ie Nash equilibrium.

A better dynamic response scheme is formed by a sequence of steps, the next step of each UAV is based on a policy applied by these competitors to previous steps, it has integrated into its process to update its policy. At the start, the first round begins with an arbitrary choice by UAV of its best response. To achieve Nash equilibrium, the following algorithm represents the Best Response learning steps that each UAV performs.

---

**Algorithm 1** Best response Algorithm

---

1: Initialize vectors $x(0) = [x_1(0), ..., x_g(0)]$ randomly;
2: **For each** $UAV_g$, $g \in \mathcal{G}$ at time instant $t$ computes:
   - $x_g(t+1) = \underset{x_g \in X_g}{\mathrm{argmax}} (U_g(x(t)))$.
3: **If** $\forall g \in \mathcal{G}$, $|x_g(t+1) - x_g(t)| < \epsilon$, then STOP.
4: **Else**, $t \leftarrow t+1$ and go to step (2)

---

Such as:
- $x$ refers to the vector price $p_c$, vector price $p_s$, vector $q_s$ or vector $\xi$.
- $X_g$ refers to the policy profile price, QoS or beconing.

## IV. NUMRICAL RESULTS

In this section, our experiments are conducted using MATLAB as a numerical simulation tool. In which the performance and efficiency of the proposed non-cooperative model are analyzed by considering two UAVs competing with each other for mobile users on the ground during a fixed period. The main objective of this work is to analyze the effect of many different parameters, such as service access rate, caching cost, QoS, and beaconing.

The analysis of funding for the implementation of the proposed approach shows that the best response algorithm

Fig. 2. Convergence of the content access fee $p_c$.



Fig. 4. Convergence of the quality of service $q_s$.



Fig. 3. Convergence of the network access fee $p_s$.



Fig. 5. Convergence of the beaconing periods $\xi$.

converges to the Nash equilibrium for fee, QoS, and beaconing periods. More importantly, we can see that with the best response algorithm, the UAV based network as expected converge quickly to the state represented by the Nash equilibrium. Figures 2, 3, 4 and 5 provide a constructive proof for the existence and uniqueness of equilibria in this setting, especially for equilibrium fee, QoS, and beaconing periods.

Both Figures 6 and 7 illustrate the fee of both network access and the content access in the equilibrium state as in function of the number of F elements. That is, as the number of elements increases, then the number of cached elements

increase as well. In addition, as the number of cached content increases, so do the caching costs. Accordingly, the UAVs will increase their fee to cover the increased cost of caching services. This conclusion is somewhat intuitive because a drone-based antenna system generates significant profits with a better pricing strategy. Furthermore, there is an effect of the number of elements F on the beaconing periods as shown in 8. As a result, the beaconing period of the proposed model increases as F increases; this is because a larger number of elements allow more content to be cached so as to easily fulfill end-users requests.

Fig. 6.  Impacts of the number of items $F$ on the content access fee.



Fig. 8.  Impacts of the number of items $F$ on the beaconing periods.



Fig. 7.  Impacts of the number of items $F$ on the network access fee.



Fig. 9.  Content access fee at different fixed cost $\mu$

In the same line, Figures 9 and 10 illustrate the access fee of both network and content for different values of caching cost, respectively. Of course, the Nash equilibrium fee is found to be minimal with a low caching cost, and from a certain value of caching cost, the network access fee and content access fee increase. One reason for this is that once the caching cost becomes expensive, the UAVs are forced to raise their fee to compensate the caching fee increases.

Figures 11 and 12 show the impact of beaconing duration on both network access fee and content access fee. As the beaconing duration increases the network access fee and

content access fee increases. The main reason is that the beaconing duration increases when the energy cost increases as well. Consequently, the UAV-based data communication in wireless sensor networks raise when the network access fee and the content access fee increase in order to compensate the increase in the energy cost.

Figures 13, 14 and 15 show the influence of backhaul bandwidth cost on network access fee, content access fee and QoS. In this respect, network access fee and content access fee increases, as the bandwidth cost gets higher. In the same manner, we note that the quality of service is decreasing with

A Novel Strategic Caching and Availability Optimization for
Wireless Unmanned Aerial Vehicle Communication Networks



Fig. 10.  Network access fee at different fixed cost $\mu$.



Fig. 12.  Content access fee versus beaconing periods $\xi$.



Fig. 11.  Network access fee versus beaconing periods



Fig. 13.  Content access fee varies with the cost $\vartheta$

the bandwidth cost. For a given bandwidth cost determined by the network provider, the UAV-based network provider is forced to decrease its fees and improve its QoS. This strategy would cause the demand for end-users requests to increase. However, as bandwidth cost increases, the UAV-based network needs to slightly increase the fee and decreases the QoS to compensate the expected increase in the bandwidth cost.

Figures 16 and 17 illustrate the influence of energy cost on the beaconing duration and QoS. Generally, as the energy cost increases, the beaconing duration and QoS decreases. Most obviously, increasing energy cost leads to lower incentives

to invest for UAVs providers in bandwidth, QoS and mobile beacon to better serve ground users.

In light of the above results of the numerical simulation, caching the popular contents in intermediate nodes has, to our knowledge, never been considered. Thus, we have demonstrated that caching-based architecture is the most cost-effective solution to increase the reliability of UAV-based communication.

## V. CONCLUSION

With the ongoing development of smart cities, using UAVs as a flying relay is expected to be the most important element

Fig. 14.   Network access fee varies with the cost $\vartheta$



Fig. 16.   The impact of the energy cost $C_b$ on the beaconing periods.



Fig. 15.   Quality of service varies with the cost $\vartheta$



Fig. 17.   The impact of the energy cost $C_q$ on the quality of service.

that will ensure the future success of reliable and robust communication systems. In this paper, we highlight the role of game theory models in designing and deployment of UAV-supported 5G network to assist and forward information to the mobile devices in full duplex. To this aim, we formulate the competition among UAVs as primary components of non-cooperative game to determine the optimal solution with the property that no single UAV can obtain a higher expected payoff. Our analysis focuses on several indicators such as beaconing time, fee to support caching of data, and QoS requirements. In this respect, we need to conduct a comprehen-

sive analysis of this competitive game among drones to prove the existence and uniqueness of Nash equilibrium that maximizes network coverage to mobile ground-based units. More importantly, a Bertrand game model with fairness concern is established, and its equilibrium fee is derived and analyzed. More practically, numerical experiments are carried out to investigate the factors, which affect the drones' strategies. When it comes to adopting a competitive strategy, theoretical analysis and simulation results show that QoS standards, fee option and competition policy have a significant influence on the expected profits. As a future work, we propose a new

REFERENCES

[1] M. E. Mkiramweni, C. Yang, J. Li, and W. Zhang, "A Survey of Game Theory in Unmanned Aerial Vehicles Communications," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 4, pp. 3386–3416, 2019, DOI: 10.1109/COMST.2019.2919613.

[2] C. Yan, L. Fu, J. Zhang, and J. Wang, "A Comprehensive Survey on UAV Communication Channel Modeling," IEEE Access, vol. 7, pp. 107 769–107 792, 2019, DOI: 10.1109/ACCESS.2019.2933173.

[3] I. A. Nemer, T. R. Sheltami, and A. S. Mahmoud, "A game theoretic approach of deployment a multiple UAVs for optimal coverage," *Transportation Research Part A: Policy and Practice*, vol. 140, pp. 215–230, Oct. 2020, DOI: 10.1016/j.tra.2020.08.004.

[4] I. U. Din, S. Hassan, M. K. Khan, M. Guizani, O. Ghazali, and A. Habbal, "Caching in information-centric networking: Strategies, challenges, and future research directions," *IEEE Communications Surveys Tutorials*, vol. 20, no. 2, pp. 1443–1474, 2018, DOI: 10.1109/COMST.2017.2787609.

[5] M. T. Nguyen, C. V. Nguyen, H. T. Do, H. T. Hua, T. A. Tran, A. D. Nguyen, G. Ala, and F. Viola, "UAV-Assisted Data Collection in Wireless Sensor Networks: A Comprehensive Survey," *Electronics*, vol. 10, no. 21, p. 2603, Oct. 2021, DOI: 10.3390/electronics10212603.

[6] Y. Yazid, I. Ez-Zazi, A. Guerrero-González, A. El Oualkadi, and M. Arioua, "UAV-Enabled Mobile Edge-Computing for IoT Based on AI: A Comprehensive Review," *Drones*, vol. 5, no. 4, p. 148, Dec. 2021, DOI: 10.3390/drones5040148.

[7] S. H. Alsamhi, O. Ma, M. S. Ansari, and F. A. Almalki, "Survey on collaborative smart drones and internet of things for improving smartness of smart cities," *IEEE Access*, vol. 7, pp. 128 125–128 152, 2019, DOI: 10.1109/ACCESS.2019.2934998.

[8] S. Koulali, E. Sabir, T. Taleb, and M. Azizi, "A green strategic activity scheduling for UAV networks: A sub-modular game perspective," *IEEE Communications Magazine*, vol. 54, no. 5, pp. 58–64, May 2016, DOI: 10.1109/MCOM.2016.7470936.

[9] A. Giagkos, E. Tuci, M. S. Wilson, and P. B. Charlesworth, "UAV flight coordination for communication networks: genetic algorithms versus game theory," *Soft Computing*, vol. 25, no. 14, pp. 9483–9503, Jul. 2021, DOI: 10.1007/s00500-021-05863-6.

[10] S. Handouf and E. Sabir, "Strategic Availability and Cost-Effective UAV-Based Flying Access Networks: S-Modular Game Analysis," *Mobile Information Systems*, vol. 2019, pp. 1–11, Jan. 2019, DOI: 10.1155/2019/4897171.

[11] S. Mehrizi, S. Chatterjee, S. Chatzinotas, and B. Ottersten, "Online Spatiotemporal Popularity Learning via Variational Bayes for Cooperative Caching," *IEEE Transactions on Communications*, vol. 68, no. 11, pp. 7068–7082, Nov. 2020, DOI: 10.1109/TCOMM.2020.3015478.

[12] D. Liu, B. Chen, C. Yang, and A. F. Molisch, "Caching at the wireless edge: design aspects, challenges, and future directions," I*EEE Communications Magazine*, vol. 54, no. 9, pp. 22–28, Sep. 2016, DOI: 10.1109/MCOM.2016.7565183.

[13] X. Xu, Y. Zeng, Y. L. Guan, and R. Zhang, "Overcoming Endurance Issue: UAV-Enabled Communications With Proactive Caching," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 6, pp. 1231–1244, Jun. 2018, DOI: 10.1109/JSAC.2018.2844979.

[14] F. Cheng, G. Gui, N. Zhao, Y. Chen, J. Tang, and H. Sari, "UAVRelaying- Assisted Secure Transmission With Caching," *IEEE Transactions on Communications*, vol. 67, no. 5, pp. 3140–3153, May 2019, DOI: 10.1109/TCOMM.2019.2895088.

[15 S. Chai and V. K. N. Lau, "Online Trajectory and Radio Resource Optimization of Cache-Enabled UAV Wireless Networks With Content and Energy Recharging," I*EEE Transactions on Signal Processing*, vol. 68, pp. 1286–1299, 2020, DOI: 10.1109/TSP.2020.2971457.

[16] H. Tran-Dinh, S. Gautam, S. Chatzinotas, and B. Ottersten, "Throughput Maximization for Wireless Communication systems with Backscatter- and Cache-assisted UAV Technology," *arXiv:2011.07955 [cs, math]*, Nov. 2020. [Online]. Available: http://arxiv.org/abs/2011.07955

[17] D. Ait Omar, M. Outanoute, M. Baslam, M. Fakir, and B. Bouikhalene, "On understanding price-QoS war for competitive market and confused consumers," *Computing*, vol. 101, no. 9, pp. 1327–1348, Sep. 2019, DOI: 10.1007/s00607-018-0642-5.

[18] D. Ait Omar, H. Garmani, F. Es-Sabery, M. El Amrani, E.-S. Azougaghe, M. Baslam et al., "Chaotic dynamics in joint price qos game with heterogeneous internet service providers," *Journal of Computer Networks and Communications*, vol. 2022, 2022, DOI: 10.1155/2022/9541887.

[19] M. El Amrani, H. Garmani, D. Ait Omar, M. Baslam, and B. Minaoui, "Analyzing the dynamic data sponsoring in the case of competing internet service providers and content providers," *Mobile Information Systems*, vol. 2021, pp. 1–16, 2021, DOI: 10.1155/2021/6629020.

[20] V. Sanchez-Aguero, F. Valera, I. Vidal, C. Tipantuña, and X. Hesselbach, "Energy-Aware Management in Multi-UAV Deployments: Modelling and Strategies," *Sensors*, vol. 20, no. 10, p. 2791, May 2020, DOI: 10.3390/s20102791.

[21] M. Mangili, F. Martignon, S. Paris, and A. Capone, "Bandwidth and Cache Leasing in Wireless Information Centric Networks: a Game Theoretic Study," *IEEE Transactions on Vehicular Technology*, pp. 679–695, 2016, DOI: 10.1109/TVT.2016.2547740.

[22] C. Fricker, P. Robert, J. Roberts, and N. Sbihi, "Impact of traffic mix on caching performance in a content-centric network," *arXiv:1202.0108[cs]*, Feb. 2012. [Online]. Available: http://arxiv.org/abs/1202.0108

[23] H. Garmani, D. A. Omar, M. E. Amrani, M. Baslam, and M. Jourhmane, "Nash bargaining and policy impact in emerging ISP-CP relationships," *International Journal of Ad Hoc and Ubiquitous Computing*, vol. 35, no. 3, pp. 117–135, 2020, DOI: 10.1504/IJAHUC.2020.110819.

[24] H. Garmani, D. Ait Omar, M. El Amrani, M. Baslam, and M. Jourhmane, "Towards a Predictive Analysis of UAV-Based Flying Base Station Decisions:," *International Journal of Business Data Communications and Networking*, vol. 16, no. 2, pp. 20–52, Jul. 2020, DOI: 10.4018/IJBDCN.2020070102.

[25] H. Garmani, D. A. Omar, M. E. Amrani, M. Baslam, and M. Jourhmane, "A game theory approach for UAV-based flying access networks," *International Journal of Networking and Virtual Organisations*, vol. 24, no. 1, p. 84, 2021, DOI: 10.1504/IJNVO.2021.111613.

[26] S. Lasaulce, M. Debbah, and E. Altman, "Methodologies for analyzing equilibria in wireless games," *IEEE Signal Processing Magazine*, vol. 26, no. 5, pp. 41–52, Sep. 2009, DOI: 10.1109/MSP.2009.933496.

[27] J. B. Rosen, "Existence and Uniqueness of Equilibrium Points for Concave N-Person Games," *Econometrica*, vol. 33, no. 3, pp. 520–534, 1965, DOI: 10.2307/1911749.

**Mohamed El Amrani** received the Ph.D. degrees from University Sultan Moulay Slimane, Morocco, in 2021. His research interests include network economics, network security, applications of game theory in wireless networks, and radio resource management.

**Hamid Garmani** is a Professor of computer science in the Faculty of Sciences and Techniques, Sultan Moulay Slimane University, Morocco. His research interests include network economics, network security, applications of game theory in wireless networks, and radio resource management.

**Driss Ait Omar** is a Professor of computer science in the Faculty of Sciences and Techniques, Sultan Moulay Slimane University, Morocco. His research interests include network economics, network security, applications of game theory in wireless networks, and radio resource management.

**Mohamed Ouaskou** is a PhD student at the Faculty of Sciences and Technology, University Sultan Moulay Slimane, Morocco. His research interests include network economics and applications of game theory in wireless networks.

**Abdelkarim Ait Temghart** is a PhD student at the Faculty of Sciences and Technology, University Sultan Moulay Slimane, Morocco. His research interests include network economics and applications of game theory in wireless networks.

**Mohamed Baslam** is a Professor of computer science in the Faculty of Sciences and Technology, Sultan Moulay Slimane University, Morocco. His current research interests include performance evaluation and optimization of networks based on game-theoretic and queuing models, applications in communication/transportation and social networks, such as wireless flexible networks, bio-inspired and self-organizing networks, and economic models of the Internet and yield management.

# Lightweight 4over6 Test-bed for Security Analysis

Ameen Al-Azzawi, and Gábor Lencse

*Abstract*—In this paper, we focus on one of the most prominent IPv6 transition technologies, namely lw4o6 (Lightweight 4over6). We emphasize the uniqueness of lw4o6 and the difference between lw4o6 and the conventional DS-Lite (Dual-Stack Lite), their topology, functionality and security vulnerabilities. We analyze the potential vulnerabilities of lw4o6 infrastructure by applying the STRIDE threat modelling technique, which stands for Spoof- ing, Tampering, Repudiation, Information Disclosure, Denial of Service, and Elevation of Privilege. Moreover, we build a test-bed for lw4o6 using Snabb, which is an open source software. We test Snabb's tunneling and binding capabilities and most importantly, port allocation per subscriber. At the end, we present multiple attacking scenarios (Denial of Service, Information Disclosure, Spoofing, etc.) against lw4o6's main routers and come up with mitigation methods for such attacks.

*Index Terms*—4in6, lw4o6, DS-Lite, encapsulation, IPv6, DoS, Snabb, PSID, Scapy.

## I. INTRODUCTION

AFTER the depletion of the public IPv4 address pool in 2011 [1], several technologies were presented by the scientific community: research institutes, industrial vendors and ISPs (Internet Service Providers). All of them aimed to fulfill one commitment, reliable communication between two remote machines that have different IP versions (IPv4 and IPv6), or the same IP version but with the infrastructure between them adopting another IP version.

We have been focusing our research on the most prominent IPv6 transition technologies and conducted a survey on them [2], where we concluded and shortlisted some of the most promising technologies. For example, the combination of NAT64 [3] and DNS64 [4] proved to be effective in certain areas. However, it did not solve the issue of IPv4-only applications. Therefore, another technology called 464XLAT [5] has been developed to tackle this problem, which has a double translation mechanism using CLAT (customer-side translator) and PLAT (provider-side translator). We have published several papers [6], [7], where we have analyzed the security threats that the 464XLAT infrastructure might face, especially CLAT and PLAT routers using the STRIDE (Spoofing, Tampering, Repudiation, Information disclosure, and Elevation of privilege) method [8].

Furthermore, we have published another article [9], where we have built a test-bed for 464XLAT, and have tested its capabilities and how it reacts to potential security threats.

A. Al-Azzawi is with the Department of Networked Systems and Services, Faculty of Electrical Engineering and Informatics, Budapest University of Technology and Economics, Budapest, Hungary. (e-mail: alazzawi@hit.bme.hu).

G. Lencse is with the Department of Networked Systems and Services, Faculty of Electrical Engineering and Informatics, Budapest University of Technology and Economics, Budapest, Hungary. (e-mail: lencse@hit.bme.hu).

The `hping3` package was used to flood the PLAT router with an excessive amount of packets. We concluded that 464XLAT is a reliable technology when it comes to IPv4-only devices communicating over IPv6 island. However, it has some security vulnerabilities that can be leveraged by an adversary, such as DoS (Denial of Service) attack.

Moreover, another technology called DS-Lite was invented to tackle the same issue of IPv4 depletion, it consists of two main routers: B4 (Basic Bridging BroadBand) and AFTR (Address Family Transition Router) [10]. We have covered the security analysis of DS-Lite in [11]. However, DS-Lite has an issue with scalability as all of the NAPT operations were executed by the AFTR router, which makes it very hard to scale its operation. Therefore, lw4o6 technology was invented to function as an improved version of DS-Lite [12]. In our current paper, the target is to conduct an analysis of the potential security vulnerabilities that the lw4o6 infrastructure may face. To that end, the STRIDE threat modeling technique [8] is used. We plan to fulfill our objectives as follows:

- Utilizing the STRIDE threat modeling technique on the Data-Flow Diagram (DFD) of a lw4o6 lw4o6 system to identify and uncover any possible security vulnerabilities.
- Constructing a reliable test-bed for the lw4o6 system using open-source software to assess its functionality and translation process.
- Conducting multiple attack scenarios on the primary routers of the lw4o6 infrastructure to evaluate their resilience, effectiveness, and potential security weaknesses.
- Ultimately, suggesting practical mitigation strategies to address and counteract such attacks.

The remainder of this paper is organized as follows. In Section II, we provide an overview of the operation of the lw4o6 infrastructure and, more importantly, the core differences between lw4o6 and DS-Lite. We emphasize the fundamental differences between the two technologies in terms of their topology, functionality, and scalability. In Section III, we talk about the previous studies that tried to build a test-bed or analyze lw4o6 transition technology. In Section IV, we describe the implementation of lw4o6 infrastructure and how it can be built. In Section V, we analyze the potential security vulnerabilities within lw4o6 using the STRIDE method after we build its DFD (Data Flow Diagram). In Section VI, we present our lw4o6 test-bed, whereas, in Section VII, we present our results and analyze them, where we also present two attacking scenarios and their mitigation methods.

In Section VIII, we conclude our paper by summarizing the significance of our test-bed and the lesson learned from its results.

Fig. 1. DS-Lite Topology.



Fig. 2. Lw4o6 Topology

## II. Lw4o6 Operation

Lw4o6 was defined in RFC-7596 [12], as an improved version of DS-Lite [10]. To understand the essence of this extended technology, we will analyze the added values that lw4o6 presents.

### A. DS-Lite Topology

Fig.1 shows the topology of DS-Lite and its encapsulation / decapsulation process. DS-Lite is based on two core functions:

- B4 (Basic Bridging BroadBand): It is responsible for encapsulating IPv4 packets into IPv6 ones and forwarding the 4in6 traffic to the AFTR (Address Family Transition Router). It also decapsulates the returning 4in6 traffic by extracting the original IPv4 packet from the payload and forwards it to the IPv4 client residing behind it [10].
- AFTR: It decapsulates the 4in6 traffic that is generated by the B4 router, then statefully translates the original private IPv4 source address into a public one [10]. Moreover, AFTR works also as an encapsulator, when it encapsulates the returning packet from the IPv4 server into IPv6 tunnel and sends it to the designated B4 router. AFTR adopts a method called per-flow state [12]. It processes every packet as a part of a flow and compares this incoming packet with its binding table by checking the softwire ID (B4 IPv6 address) then scans its binding table for translation purposes.

### B. Lw4o6 Topology

Fig. 2 shows the topology of lw4o6 with its main components. Lw4o6 works in a similar manner as DS-Lite, however, it has one main improvement, which made it more scalable than the conventional DS-Lite. The stateful translation has been moved from the centralized AFTR into the B4. This change made a big difference for the ISPs, and can be summarized by the below points:

- Lw4o6 optimizes the work for ISPs by avoiding the complicated process of stateful translation.
- The "per flow state" method is replaced by "per subscriber state", which will save a lot in terms of CPU

and memory consumption [12]. This is considered to be the most important feature of lw4o6, where every subscriber or CPE (Customer Premise Equipment) has an allocated and dedicated port set that can be used in his communication. This opens up the possibility for multiple lwB4 routers to share the same public IPv4 address in their softwire (tunnel) under the condition that they use two different port sets [12].

- The fact that no stateful translation is required on the lwAFTR side means that considerably less logging in the ISP side is actually required [12].

In other words, lw4o6 is an optimization of DS-Lite [12], as it reduces the overhead at the lwAFTR side by relocating the stateful translation to the lwB4 side.

Lw4o6 infrastructure consists of two main types of routers, lwB4 and lwAFTR:

- lwB4: works as stateful NAT44 translator and encapsulator / decapsulator.
- lwAFTR: works only as encapsulator / decapsulator.

Softwire is a mechanism that allows the encapsulation and transport of IPv4 traffic over an IPv6 network (or vice versa). Moreover, it is a tunnel ("virtual wire") that carries IPv4 traffic over an IPv6 infrastructure. In the lw4o6 case, softwire encapsulates IPv4 packets within IPv6 packets and facilitates communication between two IPv4 networks using an IPv6 network [12].

### C. Port Set Allocation

As described previously, lwAFTR allocates a specific port set for every lwB4 router using a method illustrated in RFC-7597 [13] Section 5.1, where PSID (Port Set Identifier) can be calculated as follows.

For example, to configure a softwire at the lwAFTR side, every softwire has to have the below three parameters (the numbers we used here are just examples):

- Port set size = 10 $\Rightarrow 2^{10} = 1024$ ports per set.
- PSID length = 6 $\Rightarrow$ number of port sets: $2^6 = 64$, it is also called "the sharing ratio" [13].
- PSID = 1 $\Rightarrow$ allocated ports range = [1024-2047].

To explain it in a simpler way, let us calculate the number of concurrent lwB4s (subscribers) who can share the same public IPv4 address, while having different port sets. The total number of source ports numbers is $2^{16} = 65536$, divided by 64 (number of port sets) is 1024. This number represents the size of one ports set, which can be also concluded from PSID size value ($2^{10}$). So, we have the number of sets, and the size of the sets themselves, all that is left is to select the PSID value, which will decide the exact port set to be allocated for the specific CPE (subscriber).

- PSID = 0 $\Rightarrow$ allocated ports = [0 - 1023].
- PSID = 1 $\Rightarrow$ allocated ports = [1024 - 2047].
- PSID = 8 $\Rightarrow$ allocated ports = [8192 - 9215].

In conclusion, with PSID length = 6, we can support 64 different CPEs (subscribers) with 1024 ports for each of them. Furthermore, by selecting PSID = 1, the range of allocated ports will amount to [1024 - 2047]. However, if we exclude the first set, which contains the well-known ports [0 - 1023], we end up with 63 subscribers. That means 63 subscribers have the possibility of sharing the same public IPv4 address.

## III. RELATED WORK

There is a very limited amount of research available in the field of lw4o6 IPv6 transition technology.

Ahmed Al-hamadani proposed a test environment for benchmarking lw4o6 and especially its two main components (lwB4 and lwAFTR) [14]. The author carried out an analysis for the operational requirement to build such a tester, which aimed to be the world's first RFC-8219 [15] compliant lw4o6 tester.

Omar D'yab built a test-bed for lw4o6 [16], where the author demonstrated the operation of lw4o6 with its encapsulation/decapsulation mechanism.

As for the lwB4 router's implementation, Marcel Wiget [17] has built a complete and functioning lwB4 machine, where he used several Linux commands to build the NAT44 and IPv4-in-IPv6 tunnel. In addition, this proposed lwB4 network function is isolated into its own dedicated network namespace, which gives users the flexibility and the benefit of avoiding the use of a separate VM (virtual machine) [17].

Previous trials had been carried out to build lwB4 router using OpenWrt software [18]. However, they have proven to be complicated and not reliable [17].

## IV. LW4O6 IMPLEMENTATION

### A. LwB4 Implementation

In our test-bed, we followed a similar approach to Marcel Wiget's lwB4 router implementation [17], where we used Linux commands, such as `ip -6 tunnel` and `ip route` to build the IPv4-in-IPv6 tunnel and `iptables` to implement NAPT44.

The full bash script to configure the lwB4 router is available through the "`lwB4.sh`" script in our public GitHub repository [19]. The main commands that we used to create a tunnel for encapsulation/decapsulation plus NAPT44, were as follows:

```
ip -6 tunnel add tun-lw4o6 remote 2001:db8:2::2 \
 local 2001:db8:0:1::2 mode ipip6
ip route add 192.0.2.0/24 dev tun-lw4o6 proto static
```

```
iptables -t nat -A POSTROUTING -p tcp -o tun-lw4o6 \
 -j SNAT --to 203.0.113.1
```

Below is an explanation for the IP addresses that we used in the commands above and illustrated in Fig. 3:

- 2001:db8:2::2 is the tunnel endpoint at lwAFTR side.
- 2001:db8:0:1::2 is the IPv6 address of lwB4.
- 192.0.2.0/24 is the network address of the IPv4 server.
- 203.0.113.1 is the public IPv4 address that will be used as source IP address by lwB4 when it forwards packets to lwAFTR through the softwire tunnel.

### B. LwAFTR Implementation

Several software solutions were presented to build lwAFTR routers. It was featured in VPP [20] since version v16.09. However, VPP has demonstrated complexity in its configuration, and certain modules such as lw4o6 have become outdated and lack proper maintenance from developers. In contrast, Snabb software has received better maintenance and documentation [21]. Therefore, we have decided to deploy Snabb to build our lwAFTR router. Snabb in general is a toolkit that can be used for developing network functions in user-space, which means it bypasses the kernel to process network packets [21].

While Snabb can be used on Linux systems, it does not rely on the Linux kernel networking stack. Instead, it leverages technologies like Intel's Data Plane Development Kit (DPDK) or the Solarflare OpenOnload library to directly access the NIC (Network Interface Card) and perform packet processing in user space. By bypassing the kernel, Snabb aims to achieve lower latency and higher throughput [21].

Snabb divides the machine into two separate spheres:

- Internal interface, where IPv6 packets are received and processed.
- External interface, where public IPv4 packets are forwarded to the outside world.

Snabb resides in between those interfaces and creates a binding table. The core of Snabb's configuration is a file called "lwaftr.conf", which can be used while running the "`snabb lwaftr run lwaftr.conf`" command. The full configuration script of the lwAFTR router is available through the "`lwaftr.conf`" file in our GitHub repository [19].

For a better understanding of the packet translation, encapsulation, and decapsulation process of lw4o6, we follow the packet flow in the next Subsection. It is worth mentioning that the IP scheme which we used was based on documentation IP addresses because we used Snabb in a test environment that had internet access, while we did not want to cause any sort of routing conflict.

TABLE I
LwAFTR BINDING TABLE

| Public IPv4 | PSID | PSID length | b4-IPv6 |
|---|---|---|---|
| 203.0.113.1 | 1 | 6 | 2001:db8:0:1::2 |

On the lwAFTR router, Snabb has a pre-configured binding table (see Table I), where it stores the relevant information

Fig. 3. Lw4o6 Testbed.

for every lwB4 router (every softwire), such as lwB4's public IPv4 address, IPv6 address, and the allocated port set for it.

The binding table is directly generated from the set of the configured softwires. It is never changed in response to data-plane traffic.

*C. Packet Path Through the Lw4o6 Infrastructure*

Fig. 3 shows the topology of our lw4o6 test-bed with its elements, where its operation can be summarized as below:

- The IPv4 client sends a packet with the following details:
  - Source IP address: Client IPv4 address (10.0.0.2)
  - Source port number: 5000
  - Destination IP address: IPv4 server address (192.0.2.2).
- The lwB4 receives the packet and performs the following steps:
  - NAPT44 function: the private source IPv4 address is replaced with a public IPv4 address (203.0.113.1). The source port number is replaced with an unused one from the range assigned to the subscriber. Assuming that the assigned range is 1024-2047, let the new source port number be 1050. At lwB4, the entry of the NAPT binding table is shown in Table II.
  - Encapsulation: lwB4 encapsulates the IPv4 packet into an IPv6 packet by prepending an IPv6 header to it and forwards the 4in6 packet to the lwAFTR through the softwire tunnel with the following details:
    * Source IP address: lwB4's IPv6 address: 2001:db8:0:1::2
    * Destination IP address: lwAFTR tunnel end-point IPv6 address: 2001:db8:2::2
    * Encapsulated IPv4 packet content:
      · Source IP address + port number: 203.0.113.1:1050

TABLE II
LwB4 ROUTER NAPT44 BINDING TABLE

| Private IPv4 | Source port | External IPv4 | Temporary Port | Transport Protocol |
|---|---|---|---|---|
| 10.0.0.2 | 5000 | 203.0.113.1 | 1050 | TCP |

· Destination IP address: 192.0.2.2

- The lwAFTR router receives the 4in6 packet with the same content as above. Snabb decapsulates the IPv4 packet and scans its binding table (Table I), then acts accordingly:
  - If a matching entry is found, then the IPv4 packet is forwarded to the IPv4 Internet via the external interface.
  - Otherwise the packet is dropped.
- Finally, the IPv4 packet arrives to the IPv4 server.

Packets in the reverse direction: IPv4 server ⇒ lwAFTR ⇒ lwB4 ⇒ IPv4 client, are processed as below:

- IPv4 server replies and sends the packet to lwAFTR with the following details:
  - Source IP address: 192.0.2.2
  - Destination IP address + port number: 203.0.113.1:1050
- LwAFTR receives the above packet on Snabb's external interface, then Snabb scans Table I, looking for a matching entry. Port number 1050 is a part of the port range [1024-2047], where PSID 1 refers to it explicitly. Therefore, we have a matching entry. Snabb encapsulates the IPv4 packet into an IPv6 packet using "b4-ipv6" as the IPv6 destination address and forwards the resulting 4in6 packet to the lwB4 router with the following details:
  - Source IPv6 address: 2001:db8:2::2 (tunnel end-point).
  - Destination IPv6 address: 2001:db8:0:1::2
  - Encapsulated IPv4 packet content:
    * Source IPv4 address: 192.0.2.2
    * Destination IPv4 address + port number: 203.0.113.1:1050
- LwB4 receives the above 4in6 reply packet and performs the following:
  - In case of incorrect parameters, such as port number, then the packet will be dropped immediately.
  - Decapsulates the IPv4 packet and scans its NAPT binding table (Table II) for a match. The lwB4 router then rewrites the destination IPv4 address and port

TABLE III
VULNERABILITY OF DIFFERENT DFD ELEMENTS TO DIFFERENT THREATS [1]

| | Spoofing | Tampering | Repudiation | Information Disclosure | Denial of Service | Elevation of Privilege |
|---|:---:|:---:|:---:|:---:|:---:|:---:|
| Data Flow | | ✓ | | ✓ | ✓ | |
| Data Stores | | ✓ | | ✓ | ✓ | |
| Processes | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Interactors | ✓ | | ✓ | | | |



Fig. 4. Data Flow Datagram of Lw4o6

number of the packet based on the found entry to 10.0.0.2:5000, forwarding it to 10.0.0.2.

- IPv4 client receives the packet with the below details, which concludes the packet's journey:
  - Source IPv4 address: 192.0.2.2
  - Destination IPv4 address + port number: 10.0.0.2:5000

Snabb can have multiple softwires configured in its binding table and provisioned to communicate with multiple lwB4 routers [12]. One of Snabb's (or lw4o6 in general) challenges is the liability of lwAFTR to have millions of softwires (tunnels) configured [22]. on the other hand, the capacity of the binding table could be exhausted with an excessive number of entries.

## V. SECURITY ANALYSIS

The lw4o6 technology plays a crucial role in enabling the coexistence of IPv4 and IPv6 networks during the transition phase. However, ensuring the security of lw4o6 deployments is of paramount importance to maintain the integrity and confidentiality of network communications. In this section, we conduct a comprehensive security analysis of the lw4o6 technology and evaluate its potential vulnerabilities and threats.

### A. Threat Modeling

We have selected STRIDE to be our threat modeling technique, which can be used to assess the potential security vulnerabilities of any given IT system. It was explained by A. Shostack [8]. Below is a brief overview of the STRIDE components:

- Spoofing: an attacker's impersonation of legitimate nodes and pretending to be someone else by using an innocent IP address for example and sending harmful packets [8].
- Tampering: threats related to the modification of data or configuration in the network could allow an attacker to disrupt service or steal sensitive information [8].
- Repudiation: an attacker denies the responsibility of an act such as a DNS query or money transaction [8].
- Information Disclosure: threats related to the leakage of sensitive information from the network, such as IP addresses or routing information, could allow an attacker to track or compromise nodes [8].
- DoS: an attacker floods the targeted server with an excessive amount of packets that lead to resource depletion or disruption of service in the network and prevent legitimate users from communicating with the targeted machine [8].
- Elevation of Privileges: an adversary getting access to sensitive resources by bypassing some security control protocols that lead to him gaining unauthorized access or control over the network [8].

The type of vulnerabilities depends on what is being done with the data (processing, storing, etc..) [8]. Table III shows those vulnerabilities accordingly.

### B. Applying STRIDE on lw4o6

According to [8], a DFD (Data Flow Diagram) of the examined system needs to be drawn for STRIDE to spot the potential vulnerabilities of the given system. Therefore, we

have drawn the required diagram as shown in Fig. 4, where the system has 12 vulnerable areas that can be targeted by an attacker.

### C. Attacking Possibilities

We conducted a comprehensive analysis of all potential attacking scenarios on lw4o6 by leveraging its DFD using the STRIDE method. For a detailed examination, please refer to Appendix A.

### D. Vulnerability Assessment

Upon conducting a comprehensive examination of the potential security threats associated with lw4o6, we have reached the conclusion that the exploitable threats accessible to an attacker can be succinctly summarized in Table IV. To facilitate a clearer understanding, we have classified the severity of these attacks into distinct levels, namely low, medium, high, and critical. The categorization is based on the detrimental impact inflicted upon the targeted system, as well as the intricacy involved in performing and mitigating the respective attack.

## VI. Lw4o6 Test-Bed

To build our test-bed, we used a "P" series node of NICT StarBED, Japan [23], which is a Dell PowerEdge 430 server with the following details: two 2.1GHz Intel Xeon E5-2683v4 CPUs with 16 cores each, 348 GB 2400MHz DDR4 SDRAM. We have installed a Windows 10 Pro operating system.

As shown in Fig. 3, our test-bed consists mainly of 4 machines created with Linux-based CentOS-7 Virtual machines built on top of VMware Workstation player virtualization software. Every machine has 8 GB RAM, 6 CPU cores (except lwAFTR which has 8 cores) and 20GB HDD.

As for our lwB4 and lwAFTR implementation, we already described it in details in Section IV-A and IV-B. Moreover, full configuration of lwB4 and lwAFTR (Snabb) is available in our GitHub repository [19].

## VII. Results

### A. Normal lw4o6 Process

During pinging the IPv4 server from IPv4 client side (see Fig. 3), we monitored the traffic with the `tcpdump` command on different locations as shown in Table V, where the NAT44 + encapsulation / decapsulation on lwB4 side and encapsulation / decapsulation processes on lwAFTR side are quite obvious.

### B. PSID Test

In our test-bed, we send packets from lwB4 with a specific port set [1024 - 2047]. As a result, we calibrated that on the lwAFTR side with the below details:

- PSID length : 6
- PSID : 1

Packets went through without any drop, however, when we changed the PSID value on lwaftr softwire configuration to 2, while keeping the port range on lwB4 the same, packets did not go through lwAFTR and an ICMP6 "Destination Unreachable" message was sent back to lwB4. The reason behind the packet drop is that PSID 2 refers to port set [2048 - 4095], while we kept sending packets using the old ports range [1024 - 2047].

### C. Attacking Scenarios

1) DoS Attack: As shown in Fig. 5, an attacker machine was deployed to flood the lw4o6 infrastructure with too many TCP synchronization requests and thus perform a DoS attack against lwB4 and lwAFTR.



Fig. 5. DoS Attack Against Lw4o6 Infrastructure.

The attack was executed while the IPv4 client communicating with the IPv4 server normally, and it took around 3-5 seconds for the IPv4 client to show 75% packet loss. The attack was performed using `hping3` package:

```
hping3 -S --flood -V -p 80 192.0.2.2
```

In this situation, we explored a scenario where an intruder gains entry to the client's local network and engages in harmful actions aimed at overwhelming the access point. The objective is to obstruct the client(s) from receiving responses to their requests.

Moreover, Fig. 6, shows the CPU utilization of the lwB4 machine before and after the attack, where the CPU of the machine was fully utilized within 5 seconds.



Fig. 6. CPU utilization for lwB4 Machine

2) Information Disclosure: As shown in Fig. 7, we carried out an information disclosure attack against the traffic between lwB4 and lwAFTR, where we used Scappy script to sniff the communication channel and print out the payload of the TCP/UDP packets. The attack was successful, as it printed out the content of the payload in plain text. The script can be found under the name of "`info-disclosure.py`" in our GitHub repository [19].

3) ICMP Spoofing: Since ICMP packets do not have port numbers, lw4o6 handles ICMP packets differently than it does with TCP /UDP packets, especially when it comes to packet filtering at lwB4 and lwAFTR routers. The solution

TABLE IV
SUMMARY OF THE POTENTIAL VULNERABILITIES OF LW4O6

| Attack Name | Intricacy of Performing the Attack | Intricacy of Mitigation | Attack Impact (Severity) |
|---|---|---|---|
| TCP RST Signal | Average | Average | Low |
| IP Address Spoofing | Average | Difficult | Medium |
| Packet Injection | Average | Difficult | Medium |
| Information Disclosure | Average | Easy | Medium |
| Packet's Payload Tampering | Difficult | Difficult | Medium |
| ARP Poisoning | Average | Difficult | High |
| Source Port Exhaustion | Easy | Average | High |
| TCP Session Hijack | Easy | Average | Medium |
| Network Mapping | Easy | Easy | Low |
| DoS using TCP SYN Flood | Easy | Difficult | Critical |

TABLE V
LW4O6 PACKET ENCAPSULATION/DECAPSULATION PROCESS

| Packet capture at lwB4 ens34 |
|---|
| 17:22:15.047747 IP 10.0.0.2 > 192.0.2.2: ICMP echo request, id 21142, seq 1, length 64 |
| 17:22:15.054865 IP 192.0.2.2 > 10.0.0.2: ICMP echo reply, id 21142, seq 1, length 64 |
| **Packet capture at lwB4 ens35** |
| 17:22:15.047840 IP6 2001:db8:0:1::2 > 2001:db8:2::2: IP 203.0.113.1 > 192.0.2.2: ICMP echo request, id 1026, seq 1, length 64 |
| 17:22:15.054630 IP6 2001:db8:2::2 > 2001:db8:0:1::2: IP 192.0.2.2 > 203.0.113.1: ICMP echo reply, id 1026, seq 1, length 64 1 |
| **Packet capture at lwAFTR ens35** |
| 17:22:15.050075 IP 203.0.113.1 > 192.0.2.2: ICMP echo request, id 1026, seq 1, length 64 |
| 17:22:15.051435 IP 192.0.2.2 > 203.0.113.1: ICMP echo reply, id 1026, seq 1, length 64 |

was presented in RFC-7596 Section 8.1 [12], where the lwB4 router encapsulates a port number (out of the assigned pool of ports) into the ICMP ID field. Therefore, when the ICMP packet reaches the lwAFTR router, its ICMP ID field will be inspected and treated as the source port number, where the same process of packet filtering (Section IV-C) will be applied to the packet.

The attack was made possible by a script based on a powerful Python library called Scapy [24], which is an interactive packet manipulation program. The script can be accessed under the name of "icmp-spoofer.py" in our GitHub repository [19].

As illustrated in Fig. 8, the idea behind this attack is to sniff the communication channel for any ICMP packet being forwarded from the lwB4 towards lwAFTR, then send a crafted ICMP packet to the lwAFTR.

Before sending the crafted packet, Scapy makes sure that the new packet has similar details to the original one while altering the payload's content (the transmitted data), uses the same ICMP ID (for the sake of port mapping), and then forwards it to the lwAFTR.

The packet journey of the spoofed packet follows the below path:
Attacker ⇒ lwAFTR ⇒ IPv4 Server ⇒ lwAFTR ⇒ lwB4 ⇒



Fig. 7.  Information Disclosure attack using Scapy script



Fig. 8.  ICMP / UDP packet Spoofing Attack.

IPv4 Client. As a result, the IPv4 client received a reply for a packet that it never sent.

*4) UDP Spoofing:* As shown in Fig. 8, We repeated the same attack but targeted UDP traffic. Therefore, we applied a different Python script, which can be accessed under the name of "udp-spoofer.py" in our GitHub repository [19]. The script sniffs the traffic between lwB4 and lwAFTR routers, looks for UPD traffic, and crafts new UDP packets based on the same details, especially the source port number. In addition, the newly crafted packet has the wrong payload (random text generated by the script). Eventually, the attacker machine forwards the crafted packet to the lwAFTR router. The attack was successful and the lwAFR router processed the malicious packet normally and forwarded it to the IPv4 server.

On the other hand, we reiterated the same attack with a minor adaptation. In contrast to the initial method of extracting the source port from the UDP packet, we devised a customized packet with a randomly generated source port number and subsequently directed it toward the lwAFTR. Consequently, the lwAFTR promptly discarded the packet owing to the absence of the appropriate source port. Further details and implementation of this script can be accessed in our GitHub repository under the file name "random-source-port.py" in our GitHub repository [19]

*5) Source Port Exhaustion:* Our lwB4 router was equipped with a specific port range: [1024-2047]. Therefore, we decided to exploit this vulnerability. As illustrated in Fig. 9, we used a tool called dns64perf++ that generates an excessive amount of DNS queries toward the IPv4 server [25]. DNS queries are UDP packets and they require a UDP port to be assigned for each packet. As a result, we managed to exhaust the pool ports in less than one second. Dns64perf++ tool generated 2500 packets/second, which can be found under the name of "port-exhaust.sh" in our GitHub repository [19].



Fig. 9.   Source Port exhaustion

Fig. 11 shows the last three lines of Wireshark capture on lwB4 ens35, where we managed to exhaust the source port pool [1024-2027] in less than one second. It also shows traffic stopped in less than half of a second and then resume at the 30th second. This was due to the default timeout of the UDP connection, where ports are re-assignable after 30 seconds.

It is worth mentioning that DNS64perf++ was designed to be used as a measurement tool, however, we used it as an attacking method.



Fig. 10.   TCP Session Hijacking Attack

*6) TCP Reset Signal:* The idea behind the attack is to send a TCP RST signal to the IPv4 server (via the lwAFTR router) and terminate an existing TCP connection with the IPv4 client. We achieved our goal by writing a Python script that sniffs the traffic between lwB4 and lwAFTR routers, searches for TCP packets with SYN and ACK flags, and crafts new TCP RST packet accordingly. The crafted TCP RST packet was forwarded by the attacker to the lwAFTR router, which forwarded the packet again to the IPv4 server. The crafted packet forced the IPv4 server to terminate the TCP connection with the IPv4 server (until the next TCP Sync packet comes again from the IPv4 client).

The script can be found under the name of "lw-tcp-reset.py" in our GitHub repository [19].

*7) TCP Session Hijacking:* TCP session hijacking attack is a cyberattack where an unauthorized party intercepts and takes control of an established TCP connection between two communicating entities, potentially gaining unauthorized access and control over the communication or data exchange [26].

We conducted this attack by employing the Scapy software to intercept communication between lwB4 and lwAFTR machines, as depicted in Figure 10. The attack is initiated when the "attacker-1" machine detects a TCP ACK packet moving from the IPv4 client to the IPv4 server. Upon identification of this packet, the attacker extracts pertinent information such as IP addresses, port numbers, sequence, and acknowledgment numbers from the exchanged TCP packets and sends crafted packets to the IPv4 server accordingly.

Subsequently, the "attacker-1" machine fabricates a TCP packet with the "RST" flag and relays it to the IPv4 client via the lwB4 machine in order to deceive the IPv4 client with a fake TCP session abortion signal. Furthermore, the attack script continuously monitors the communication channel, responding to relevant TCP packets and adapting responses with accurate information and appropriate flags.

Concurrently, the script establishes an SSH connection from the "attacker-1" to the "attacker-2" machine and initiates a DoS attack against the lwB4 router with the goal of exhausting its CPU resources. This action aims to obstruct any future communication between the IPv4 client and the IPv4 server via the lwB4 router. Consequently, the "attacker-1" machine gains the ability to communicate with the IPv4 server and hijacks the ongoing TCP session.

| No. | Time | Src. IP | Dst. IP | Protocol | Src. port | Dst. port | Info. |
|-----|------|---------|---------|----------|-----------|-----------|-------|
| 1028 | 0.561305919 | 203.0.113.1 | 192.0.2.2 | DNS | 2045 | 53 Standard query 0x03fd AAAA 000-000-003· |
| 1029 | 0.562221123 | 203.0.113.1 | 192.0.2.2 | DNS | 2046 | 53 Standard query 0x03fe AAAA 000-000-003· |
| 1030 | 0.562592491 | 203.0.113.1 | 192.0.2.2 | DNS | 2047 | 53 Standard query 0x03ff AAAA 000-000-003· |
| 1046 | 30.010075937 | 203.0.113.1 | 192.0.2.2 | DNS | 1024 | 53 Standard query 0xd524 AAAA 000-000-213· |

Fig. 11. Wireshark Capture on lwB4 ens35.

Notably, the script is also designed to send meticulously crafted TCP packets containing the "`PSH`" and "`ACK`" flags (with falsified payloads) to the IPv4 server. This process is iterated each time the IPv4 server responds with a TCP packet bearing an "`ACK`" flag. The script can be found under the name "`tcp-session-hijack.py`" in our GitHub repository [19].

It's important to note that the use of the "`PSH`" and "`ACK`" flags is not mandatory in all TCP connections. By utilizing it, the server can ensure that data is delivered to the application as soon as possible, reducing latency, and improving the user experience. Results of the attack were made public under this path: "`files/tcp-session-hijack-results.txt`" in our GitHub repository [19], where it shows the full chain of TCP communication between the attacker and the IPv4 server via the lwAFTR machine.

*8) Network Mapping:* In order to expose the topology of lw4o6, we conducted an experiment with a simple attacking command:

```
traceroute 192.0.2.2
```

The result was as follows:

```
30 hops max, 60 byte packets
1 10.0.0.1 (10.0.0.1) 1.154 ms  0.739 ms  1.911 ms
2  192.0.2.2 (192.0.2.2)  36.101 ms  36.302 ms \
   36.627 ms
```

Such a result tells the attacker that there are two hops till he can reach the target. It also can help the attacker to map the network topology and reveals the network infrastructure, routers, and servers in between the source and target. This knowledge can aid in identifying potential points of entry or weak links in the network.

*D. Attacks summary*

In conclusion, by employing the STRIDE method to analyze lw4o6 security, our investigation has successfully established a link between STRIDE and actual testbed attack scenarios. We have identified and validated the following attack possibilities:

- A DoS attack (using TCP SYN flood) directed at the lwB4 machine, detailed in Appendix A4: (v).
- An Information Disclosure attack against the traffic between lwB4 and LwAFTR routers, detailed in Appendix A6: (ii).
- Two different types of Spoofing attacks targeting the lwB4 IPv6 address, detailed in Appendix A4: (i).
- A source port exhaustion attack against the lwB4 router, detailed in Appendix A2: (iii).
- A TCP RST signal attack against the lwAFTR router / IPv4 server, detailed in Appendix A6: (i).

- TCP session hijacking attack against the lwAFTR router, detailed in Appendix A6: (i).
- Network Mapping attack against the lwB4 router, detailed in Appendix A4: (iv).

*E. Mitigation Methods*

*1) DoS Attack Mitigation:* We mitigated the attack by deploying `iptables` rules on the lwB4 machine to perform a rate-limiting mechanism and setting a filter with two rules, one that allows packets to be forwarded with certain limits (10 p/s for instance), while the second rule drops everything else.

```
iptables -A FORWARD -p tcp --syn -m limit \
 --limit 10/s -j ACCEPT
iptables -A FORWARD -p tcp --syn -j DROP
```

*2) Information Disclosure Mitigation:* We have successfully mitigated the information disclosure attack, which we presented in Section VII-C2. We achieved that by encrypting the payload of TCP/UDP packets generated by the IPv4 client using the Python Fernet module. The script can be found under the name of "`payload-encryptor.py`" in our GitHub repository [19]. Therefore, after running the same attacking script, the actual payload of the packet was not visible, only its encrypted value.

*3) ICMP Spoofing Mitigation:* An advanced packet crafting software packages, such as Scapy [24], can create a very realistic crafted packet with almost the exact anticipated values (IP addresses, port numbers, Packet sequence, and even MAC addresses). Therefore, a sophisticated tool such as SNORT [27], which functions as IDS (Intrusion Detection System) and IPS (Intrusion Prevention System) as well, would be required. The tool performs a deep inspection mechanism, where it detects the suspicious (spoofed) packet and drops it eventually. SNORT proved to be complicated to configure. Therefore, we omitted it as it is out of our scope.

*4) UDP Spoofing Mitigation:* Please refer to the above Subsection (VII-E3).

*5) Source Port Exhaustion Mitigation:* To counter the attack, we implemented rate limiting for DNS queries, ensuring that only 100 packets per second are allowed. We achieved this by configuring "`iptables`" rules to drop incoming DNS query packets by default while permitting the specified rate (100 packets per second).

As a result, the pool of the allocated ports within the lwB4 router could not be exhausted with such a low rate of DNS queries. The complete script can be found under the name of "`exhaust-mitigate.sh`" in our GitHub repository [19].

*6) TCP Reset Signal Mitigation:* The mitigation proved to be very complicated, an IDS / IPS device is required for deep inspection, please refer to the above Sub-section (VII-E3).

*7) TCP Session Hijack:* Please refer to the above Sub-section (VII-E3).

*8) Network Mapping:* There are several possible mitigation methods for such an attack. For example, disabling ICMP echo replies at the network perimeter or on specific routers, will prevent outsiders from using "`traceroute`" to discover the internal network structure. Moreover, network segmentation strategy can help protect against network mapping attacks and enhance overall network security. Network segmentation involves dividing a large network into smaller, isolated segments or sub-networks. Each segment is logically separated from the others and may have its own security policies, access controls, and communication rules [28].

## VIII. CONCLUSION

Lw4o6 proved to be an optimization over DS-Lite, that reduces the overhead from lwAFTR. It also implements IPv4aaS (IPv4 as a service) in IPv6 only environment. The fact that it is stateless in the center of the network, made it possible to scale its routers capabilities, especially lwAFTR. Lw4o6 can be built using open-source software packages such as VPP, Snabb, etc. Moreover, Snabb proved to be reliable software to build lwAFTR over Linux-based virtual machines. Our test-bed was a success, where we performed the normal packet NAT44 operation + packet encapsulation/decapsulation smoothly. We managed to find some vulnerabilities using our attacking scenarios such as DoS, spoofing, Tampering and Information Disclosure, where we implemented and proposed some mitigation methods against them.

## ACKNOWLEDGMENT

## REFERENCES

[1] G. Lencse, Y. Kadobayashi, Methodology for the identification of potential security issues of different IPv6 transition technologies: Threat analysis of DNS64 and stateful NAT64, Computers & Security 77 (2018) 397–411. DOI: 10.1016/j.cose.2018.04.012.

[2] A. Al-Azzawi, Towards the security analysis of the five most prominent IPv4aas technologies, Acta Technica Jaurinensis 13 (2) (2020) 85–98. DOI: 10.14513/actatechjaur.v13.n2.530.

[3] M. Bagnulo, P. Matthews, I. van Beijnum, Stateful NAT64: Network Address and Protocol Translation from IPv6 Clients to IPv4 Servers, RFC 6146, IETF (2011). DOI: 10.17487/RFC6146.

[4] M. Bagnulo, A. Sullivan, P. Matthews, I. Van Beijnum, DNS64: DNS Extensions for Network Address Translation from IPv6 Clients to IPv4 Servers, RFC 6147, IETF (2011). DOI: 10.17487/RFC6147.

[5] M. Mawatari, M. Kawashima, C. Byrne, 464XLAT: Combination of Stateful and Stateless Translation, IETF RFC 6877, IETF (2013). DOI: 10.17487/RFC6877.

[6] A. Al-Azzawi, G. Lencse, Towards the Identification of the Possible Security Issues of the 464XLAT IPv6 Transition Technology., in: TSP, 2020, pp. 439–444. DOI: 10.1109/TSP49548.2020.9163487.

[7] A. AL-Azzawi, G. Lencse, Testbed for the Security Analysis of the 464XLAT IPv6 Transition Technology in a Virtual Environment, in: 2021 44th International Conference on Telecommunications and Signal Processing (TSP), IEEE, 2021, pp. 5–9. DOI: 10.1109/TSP52935.2021.9522598.

[8] A. Shostack, Threat modeling: Designing for security, John Wiley & Sons, 2014.

[9] A. Al-Azzawi, G. Lencse, Identification of the Possible Security Issues of the 464XLAT IPv6 Transition Technology, Infocommunications Journal 13 (4) (2021) 10–18. DOI: 10.36244/ICJ.2021.4.2.

[10] A. Durand, R. Droms, J. Woodyatt, Y. Lee, Dual-stack lite broadband deployments following IPv4 exhaustion, Rfc 6333, IETF (2011). DOI: 10.17487/RFC6333.

[11] A. Al-Azzawi, G. Lencse, Analysis of the Security Challenges Facing the DS-Lite IPv6 Transition Technology, Electronics 12 (10) (2023). DOI: 10.3390/electronics12102335. URL https://www.mdpi.com/2079-9292/12/10/2335

[12] Y. Cui, Q. Sun, M. Boucadair, T. Tsou, Y. Lee, I. Farrer, Lightweight 4over6: An extension to the dual-stack lite architecture, Rfc 7596, IETF (2015). DOI: 10.17487/RFC7596.

[13] O. Troan, W. Dec, X. Li, C. Bao, S. Matsushima, T. Murakami, T. Taylor, Mapping of address and port with encapsulation (MAP-E), Rfc 7597, IETF (2015). DOI: 10.17487/RFC7597.

[14] A. Al-hamadani, G. Lencse, Design of a software tester for benchmarking lightweight 4over6 devices, in: 2021 44th International Conference on Telecommunications and Signal Processing (TSP), 2021, pp. 157–161. DOI: 10.1109/TSP52935.2021.9522607.

[15] M. Georgescu, L. Pislaru, G. Lencse, Benchmarking methodology for IPv6 transition technologies, Rfc 8219, IETF (2017). DOI: 10.17487/RFC8219.

[16] O. D'yab, G. Lencse, Testbed for the Comparative Analysis of DS-Lite and Lightweight 4over6 IPv6 Transition Technologies, in: 2022 45th International Conference on Telecommunications and Signal Processing (TSP), IEEE, 2022, pp. 371–376. DOI: 10.1109/TSP55681.2022.9851309.

[17] D. P. García, The B4 network function, accessed: 2022-08-15 (2018). URL https://blogs.igalia.com/dpino/2018/02/15/the-b4-network-function/

[18] OpenWrt Software, http://openwrt.org (2004).

[19] A. Al-Azzawi, Lightweight 4 over 6 test-bed (9 2022). URL https://github.com/ameen-mcmxc/lw4o6-automation

[20] FDio, VPP software (2016). URL https://github.com/FDio/vpp

[21] D. P. Garcia, Snabb explained in less than 10 minutes, accessed: 2022-09-04 (2017). URL http://blogs.igalia.com/dpino/2017/11/13/snabb-network-toolkit/

[22] D. Garcia, Lightweight 4-over-6: a compelling IPv4+IPv6 architecture, accessed: 2023-04-24 (2017). URL https://www.igalia.com/project/lw4o6

[23] Making a synthesis emulation in IOT ERA possible Starbed5 Project. StarBED5 Project website, accessed: 2022-09-05 (2023). URL https://starbed.nict.go.jp/en/equipment/

[24] P. Biondi, Scapy, Packet Manipulation Program, accessed: 2023-01-02 (2003). URL https://scapy.net/index

[25] B. Dániel, DNS64perf++ Measurement Tool (2016). URL https://github.com/bakaid/dns64perfpp

[26] Z. Qian, Z. M. Mao, Off-path TCP Sequence Number Inference Attack - How Firewall Middleboxes Reduce Security, in: IEEE Symposium on Security and Privacy, 2012, pp. 347–361. DOI: 10.1109/SP.2012.29.

[27] B. Caswell, J. Beale, Snort 2.1 intrusion detection, Elsevier, 2004.
[28] N. Mhaskar, M. Alabbad, R. Khedri, A Formal Approach to Network Segmentation, Computers & Security 103 (2021) 102162. DOI: 10.1016/j.cose.2020.102162.

[29] S. Naval, V. Laxmi, M. Rajarajan, M. S. Gaur, M. Conti, Employing program semantics for malware detection, IEEE Transactions on Information Forensics and Security 10 (12) (2015) 2591–2604.

[30] C. L. Abad, R. I. Bonilla, An Analysis on the Schemes for Detecting and Preventing ARP Cache Poisoning Attacks, in: 27th International Conference on Distributed Computing Systems Workshops (ICDCSW'07), IEEE, 2007, pp. 60–60. DOI: 10.1109/ICDCSW.2007.19.

[31] S. Deng, X. Gao, Z. Lu, X. Gao, Packet Injection Attack and Its Defense in Software-Defined Networks, IEEE Transactions on Information Forensics and Security 13 (3) (2018) 695–705. DOI: 10.1109/TIFS.2017.2765506.

**Ameen Al-Azzawi** received his MSc in Communication Engineering from Northumbria University, Newcastle, England in 2014. He is now a PhD student at the Budapest University of Technology and Economics, Budapest, Hungary. He has been working full-time on his research at MediaNets Laboratory in the Department of Networked Systems and Services since September 2019. His research focus is on IPv6 transition technologies and their security analysis.

**Gábor Lencse** received his MSc and PhD in computer science from the Budapest University of Technology and Economics, Budapest, Hungary in 1994 and 2001, respectively.

He has been working full-time for the Department of Telecommunications, Széchenyi István University, Győr, Hungary since 1997. Now, he is a Professor. He has been working part-time for the Department of Networked Systems and Services, Budapest University of Technology and Economics, Budapest, Hungary as a Senior Research Fellow since 2005.

His research interests include the performance and security analysis of IPv6 transition technologies.

## APPENDIX

### A. Lw4o6 Attacking possibilities using STRIDE

#### 1) IPv4 Client:

(i) Spoofing: an attack spoofs the source IP address of the IPv4 client and misuses the possession of such source IP address by sending harmful packets towards the lwB4 router [8].

(ii) Repudiation: an attacker denies the responsibility of doing any sort of activity such as sending a malicious packet towards the lwB4 router or a packet with a spoofed IP address [8].

#### 2) Data Flow from IPv4 Client towards lwB4 Router:

(i) Tampering: an attacker intercepts and alters the data being sent from the client to the lwB4 router such as changing the client's original sent information or manipulating an order placed on an e-commerce website [8].

(ii) Information Disclosure: an attacker intercepts and views sensitive information being sent by the client to the lwB4 router, such as viewing a client's credit card information or login credentials [8].

(iii) DoS: an attacker floods the lwB4 router with a large number of requests from the IPv4 client side (or multiple clients), causing the router to become overwhelmed and unable to process legitimate requests. This can cause the router to crash or become unresponsive and deny service to legitimate clients [8]. Moreover, an attacker might send an excessive amount of DNS queries in order to exhaust the allocated pool of ports for the lwB4 router.

#### 3) Data Flow from lwB4 Router to IPv4 Client:

(i) Tampering: an attacker alters data sent from the lwB4 router to the client in order to disrupt or gain unauthorized access to the IPv4 client's system. For example, an attacker might modify a software update file sent from a server to the IPv4 client in order to include malware [29]

(ii) Information Disclosure: an attacker intercepts or otherwise gains access to sensitive information sent from a server to the IPv4 client. For example, an attacker might use a man-in-the-middle attack to intercept and read login credentials sent from a server (via lwB4 router) to the IPv4 client in clear text [8].

(iii) DoS: an attacker floods the IPv4 client with an excessive amount of requests from the lwB4 router side, causing the client to become unavailable to legitimate incoming requests. For example, an attacker might launch a distributed denial of service (DDoS) attack against the client in order to disrupt its access to a specific website [8].

#### 4) The lwB4 Router:

(i) Spoofing: an attacker spoofs the source IP address of the lwB4 router (impersonates it) and initiates communication with neighbouring (or remote) devices, while sending all sorts of malicious packets, that could harm the reputation of the organization that operates the router itself [8]. Moreover, another attack scenario is possible such as the potential for an Address Resolution Protocol (ARP) cache poisoning attack. The attack involves the interception of network traffic between the lwB4 router and the lwAFTR router, enabling an attacker to exploit the situation. By sending deceptive ARP messages to the lwB4 router, the attacker can assume the identity of the lwAFTR router. Consequently, the lwB4 router updates its ARP cache with the attacker's Media Access Control (MAC) address, erroneously associating it with the lwAFTR router. Once the ARP cache has been successfully poisoned, the attacker gains the ability to intercept, manipulate, or extract sensitive information from the network traffic. This form of attack poses a significant threat, capable of compromising the security of the system [30].

(ii) Tampering: an adversary might alter the actual data that is being processed or stored within the router such as IP addresses, port numbers, TTL (Time To Live) values, etc. It may lead to redirecting the packet to a malicious server and prevent the legitimate recipient of the packet from getting a response to his request [8].

(iii) Repudiation: after spoofing the source IP address of the lwB4 router and sending harmful packets, the attacker will be able to deny the fact that he was behind those suspicious packets [8].

(iv) Information Disclosure: an attacker having unauthorized access to confidential data within the lwB4 router such as packet's payload or the routing table of the lwB4 router itself [8]. Another type of information disclosure attack can be performed by an attacker gaining access to the network topology and the number of hops to reach a specific target.

(v) DoS: an attacker might target the lwB4 router with the well-known DoS attack. This can be done either from the IPv4 client side or from the lwAFTR side, where a huge amount of useless packets can be sent to the lwB4 router with the aim to overwhelm its processing power [8], please refer to Appendix A3.(iii).

(vi) Elevation of Privileges: an adversary bypasses the authority matrix within the organization to gain access (such as read/write permission) and therefore pushes a destructive configuration to the lwB4 router, which affects the whole network [8].

*5) The NAPT44 Binding Table of lwB4 Router:*

(i) Tampering: an attacker could potentially tamper with the NAPT44 binding table (see Table II) by altering the information stored in it, such as the IP addresses, ports, or protocol information. This causes the lwB4 to misroute or block legitimate traffic, leading to a communication breakdown [8].

(ii) DoS: an attacker could launch a Denial of Service attack on the NAPT44 binding table by overloading it with a large number of fake or malformed connection entries. This could cause the table to become full, preventing the lwB4 router from tracking legitimate connections and resulting in a communication breakdown [8].

*6) Data Flow from lwB4 to lwAFTR:*

(i) Tampering: an attacker has the capability to introduce a malevolent packet into the network traffic, leading to detrimental consequences for the lwAFTR. This act can be classified as a packet injection attack [31]. As an example, the malicious packet could manifest as a TCP RST packet, which promptly terminates an existing TCP connection. Moreover, the attacker might hijack the TCP session by eavesdropping on the communication between the lwB4 and the lwAFTR, then start communicating with the IPv4 server via the lwAFTR.

(ii) Information Disclosure: an attacker intercepts and views sensitive information, such as the TCP/UDP packet's payload, please refer to Appendix A2: (ii).

(iii) DoS: an attacker floods the lwAFTR router with useless packets to overwhelm it, please refer to Appendix A2: (iii).

*7) Data Flow from lwAFTR to lwB4:*

(i) Tampering: an attacker performs a man-in-the-middle attack and alters the packet's details such as the destination IP address, which will re-direct the packet to a potentially malicious server and deprive the legitimate lwB4 router of the response it was anticipating.

(ii) Information Disclosure: as the attacker performs a man-in-the-middle attack, he also exposes the content of the sent data, which is a data confidentiality breach.

(iii) DoS: an attacker floods the lwB4 router with an excessive amount of packets in order to overwhelm its computation power, please refer to Appendix A2: (iii).

*8) The lwAFTR Router:*

(i) Spoofing: an attacker might impersonate the lwAFTR router and initiate a communication with the IPv4 Server or the lwB4, which will make all of those network elements liable to the risk of exchanging sensitive data with an adversary [8].

(ii) Tampering: an attacker alters the content of the sensitive data within the lwAFTR router, such as the provisioned PSID value assigned for a specific lwB4 router, please refer to Appendix A4: (ii).

(iii) Repudiation: the packet sender hides his own identity, please refer to Appendix A4: (iii).

(iv) Information Disclosure: an attacker getting access to confidential data inside the lwAFTR router, such as packet's payload [8], please refer to Appendix A4: (iv).

(v) DoS: an attacker can exhaust the computation power of the lwAFTR router by sending too many useless packets, which will prevent it from being able to process any incoming packets from the lwB4 router.

(vi) Elevation of Privileges: an attacker getting high privilege access to the lwAFTR router such as admin permission [8]. Such attacks happen most of the time due to inside job [1].

*9) LwAFTR Binding Table:* This table stores all of softwires (tunnels) configured into the lwAFTR, see Table I.

(i) Tampering: an attacker alters the content of an entry within the table such as PSID value or lwB4's public IPv4 address. Such changes will pave the way for a malicious spoofed packet to be processed by the lwAFTR, while dropping the legitimate requests.

*10) Data Flow from lwAFTR Router towards IPv4 Server:*

(i) Tampering: please refer to Appendix A3: (i)

(ii) Information Disclosure: please refer to Appendix A3: (ii)

(iii) DoS: please refer to Appendix A3: (iii)

*11) Data Flow from IPv4 Server towards lwAFTR Router:*

(i) Tampering: please refer to Appendix A2: (i).

(ii) Information Disclosure: please refer to Appendix A2: (ii).

(iii) DoS: please refer to Appendix A2: (iii).

*12) IPv4 Server:*

(i) Spoofing: please refer to Appendix A1: (i).

(ii) Repudiation: please refer to Appendix A1: (ii).

# Holistic attack methods against power systems using the IEC 60870-5-104 protocol

János Csatár, Péter György, and Tamás Holczer

*Abstract*—IEC 60870-5-104 is a widely used protocol for telecontrol in European power systems. However, security was not a design goal when it was originally published: This protocol lacks built-in security features such as encryption, integrity protection, or authentication. In this paper, we describe novel types of attacks against the protocol in a holistic way. Therefore, we also enumerate the possible entry points of the threat actors and demonstrate a new technique, where the malicious actor can precisely target the attack. These methods are demonstrated both on simulated environment and actual devices and compared with already published methods.

*Index Terms*—IEC 60870-5-104, Attack, Security, Power system

## I. INTRODUCTION

As power systems became more interconnected and increasingly complex, the operation of it started to rely more and more on automation, remote sensing and control in favor of efficiency and reliability. This escalated to a point where the usage of information and communication technologies (ICT) became an essential part of power system operation; today's large power systems are already unable to function properly without ICT. Thus, it is a cyber-physical system where the process of physical power transmission and distribution is closely intertwined with communication. Since power systems have their own special needs, some industry specific communication protocols exist. One of these is IEC 60870-5 series standard, and its companion standard IEC 60870-5-104 [1] for TCP/IP based communication (abbreviated as IEC 104 further on). It was designed to transmit timestamped counters, measurement values, status signals and control commands or set-points and is widely used in some countries for SCADA control centers handling the operation of the power system; both on transmission and distribution levels. Cyber security was not a priority at the early days of power system digitization - and so core IEC 104 standard is lacking security measures. A security extension is given for 60870-5 based protocols in IEC 62351 [2] series standard incorporating - among others - standard TCP/IP security measures. Subsequently, supplemental information on implementation is now part of the IEC 60870-5

series (60870-5-7). However, two main challenges hinder the universal use of the more secure version.

- First, it requires hardware level support.
  - Several legacy equipment and software component exist due to the long life cycle that is usual in power systems.
  - Legacy support means coexistence of insecure and secure versions of IEC 104.
- The second challenging factor is increased complexity regarding a number of factors.
  - Installation and maintenance of more secure systems means increased effort, for example: managing role based access, certificate authorities, keys, etc.
  - Monitoring and validating traffic on the network becomes less straightforward and more time/resource consuming.
  - Debugging and investigation of malfunction becomes a multi level task.
  - The communication itself becomes more resource intensive - it involves overhead and so introduces latency and heavier network load.

But even in a fully IEC 62351 compliant system, security measures can still have vulnerabilities, which exposes the *core IEC 104* communication flow. The newer, paradigm changing IEC 61850 series standard is currently still more focused on substations than control centers; even some guideline exist in the IEC 61850 series how to map objects to other lower level communication protocols for usage beyond substations, like IEC 104.

Summing the above, IEC 104 will remain in operation for the coming years. It is mainly used between remote terminal units (RTU) and supervisory control and data acquisition (SCADA) systems. Since it is based on TCP/IP, an ICT infrastructure is needed, that is generally multipurpose and handles other types of traffic (e.g. remote access, configuration, network monitoring, etc.). At the same time, the SCADA system itself can be even geographically spread out; it is usual that the backup system is located at a different site.

The last few years showed that well resourced and sophisticated, power system targeted attacks can happen, like the two Ukrainian attacks in 2015 and 2016 [3], [4]. The case with SolarWinds [5], or log4j [6] pointed out that no system or component of a system can be considered completely secure as a vulnerability may be introduced by third party components. Even the communication provider of critical services can be attacked as we have seen recently in Portugal [7]. [8] shows a comprehensive survey of the topic. Compromising

even part of a power system's IEC 104 communication could give substantial control to an attacker. All this makes the cybersecurity investigation of IEC 104 based communication flow a current research topic.

The focus of the research was put on the creation of a comprehensive attack scenario against a system using IEC 60870-5-104. To do so we:

1) enumerated the possible attack surfaces, where the attack can begin,
2) created a new, sophisticated way of packet injection into a live connection,
3) proposed how the communication nodes could be identified with open-source intelligence (OSINT) and other data by the adversary.

Furthermore, we demonstrated the effective use of the proposed packet injection attack on both simulated and real equipment.

The remainder of the paper is structured in the following way: We give a short overview of IEC 104 for the reader's convenience in Section II and discuss the related work afterwards. In Section IV we give an overview of possible entry points of an attack in case of power systems. The known and proposed new attacks against IEC 104 are described in details in Section V. The problem and a promising solution of precisely targeting such an attack is analysed in Section VI. Finally we conclude our paper in Section VII.

## II. OVERVIEW OF THE IEC 60870-5-104 PROTOCOL

The IEC 60870-5-104 (IEC 104) is part of the IEC-60870 protocol family and is widely used for power systems in Europe (DNP3 is more widely used for the same purpose in North America). IEC 104 defines network access using TCP for IEC 60870-5-101.

The IEC 104 messages (Application Protocol Data Unit, APDU) are divided into two parts:

- APCI (Application Protocol Control Information) - Message header
- ASDU (Application Service Data Unit) - Message body

An APDU can contain a single APCI or an APCI and an ASDU. Generally, the length of the APCI is 6 bytes. In the remainder of this section, we will introduce the basics (from our point of view) of IEC 104 APDUs. These basics are necessary to understand why and how does some protocol-specific attacks operate. An in-depth analysis of IEC 104 was published in [9], which describes the protocol in more detail from every aspect.

### A. Application Protocol Control Information

Each APCI (Application Protocol Control Information) starts with a start byte with value 0x68, followed by the 8-bit length of APDU (Application Protocol Data Unit) and four 8-bit control fields(CF).

The frame format is determined by the two last bits of the first control field. The standard defines three frame formats, I-format, U-format, and S-format. The S and I format stores the sequence numbers about the messages sent. If this counter is invalid, then the connection is terminated. Figure 1 shows the structure of different frames.



Fig. 1. The frame formats of APCI [9]

### B. Application Service Data Unit

The ASDU contains two main sections:

- Data Unit Identifier
- Payload

The ASDU starts with a type of identification, which specifies the command. It also stores the originator's address, the addressee's address (ASDU address), the information object address, and the information element (value).

The ASDU and the information object address (IOA) are used as selectors. The ASDU address is a two-byte long address field that specifies the recipient device. The IOA is a three-byte long address, and it determines which objects are affected by the command.

The operators need to know both the ASDU address and the IOA to send messages. The ASDU addresses and the IOAs of the devices are not publicly accessible. The values are assigned by the architects of the power grid and may differ for each power grid.

As Figure 2 shows, the ASDU address and the IOA are used to identify which element is affected by the message. Even though there are multiple devices with the same IOA, only one gets the message.

There are 127 different messages used in IEC 104, both monitoring and control are covered. When a server receives a command from the client, it responds based on the requests (acknowledgment for control commands, values of sensors for monitor commands). Figure 3 shows the structure of the ASDU.

Fig. 2. IEC 60870-5-104 message identifiers.



Fig. 3. The structure of ASDU [9].

It can be seen from this brief introduction, that this protocol does not provide any cryptographic protection for the messages being sent. This property enables different kinds of attacks described in Section V.

## III. RELATED WORK

The IEC-60870 series is a widely researched and used protocol family, especially the IEC-60870-5-104 (IEC 104) protocol. Research usually focuses on how the protocol can be used efficiently, but the security aspects of the protocol are also studied. Most of the security-oriented research focuses on anomaly and attack detection in IEC 104 networks [10], [11], [12], [13], [14].

These detection mechanisms usually rely on deep packet inspection or data and time correlation. [14] compares different learning algorithms for anomaly detection and concludes that no unsupervised algorithm works perfectly alone. [13] focuses on detection of denial of service attacks based on machine learning. The paper shows by experiments that decision trees are suitable for DoS detection. In [10] shows that in some cases the Rule Learner classifier algorithm works well. [11] shows that timing anomalies in IEC 104 traffic can be detected efficiently if the anomaly is persistent or at least one hour long. Shorter timing anomalies are hard to detect. [12] proposes a traditional rule based algorithm that can efficiently detect known attacks. Recent reviews [15], [16] compares the different machine learning methods for anomaly detection in SCADA systems. Some of these methods can be used to detect attacks against the IEC 104 protocol. Apart from the detection, there is also a paper about hardening of the protocol [17].

Not just the defense side of the topic was researched, but multiple papers were published about attack vectors, and flaws in the protocol [18], [19], [20] and against systems that use IEC 104 [21], [22], [23]. Also, some exciting attack scenarios against power grid control infrastructure contain steps against the IEC 104 protocol [24]. During our research, we used the attacks described in these papers as inspiration. Furthermore, we tried to make them more sophisticated, and we even came up with new ideas that can be used in an attack scenario. In the following, we introduce the most relevant papers and position our current work.

Multiple flaws and attack vectors are described in [18] including unauthorized access, denial of service, and man-in-the-middle (MitM) attacks. As mentioned in the analysis of the IEC 104 protocol, it does not provide authentication. Therefore somebody with communication opportunities with the server can also use the services provided by IEC 104 and act as a legitimate IEC 104 client. The lack of authentication can result in devastating attacks because even with little technical knowledge, an attacker could cause severe harm to the network.

The authors of [18] tried two different attacks to cause a denial of service on the IEC 104 server. They tried flooding the server with commands which can be used to confuse the operator. The other attack was a TCP SYN DoS attack. The idea is the same as in a typical SYN DoS attack, they sent SYN packets to the IEC 104 server, and it responded with SYN-ACK, which the initiator of the communication did not handle. Using this attack, they could significantly increase the hardware and network usage of the device. Because the devices used in an IEC 104 network usually have minimal hardware resources, this attack can be efficiently used to cause a denial of service.

The authors also simulated a MitM attack, where they achieved the MitM position using ARP poisoning. They used ettercap[1] for carrying out ARP poisoning. Furthermore, ettercap filters were used to isolate and drop every IEC 104 packet, therefore terminating the communication of the legitimate parties. The authors of [18] used the OpenMUC implementation of the protocol to demonstrate these attacks.

---

[1] https://www.ettercap-project.org/

A similar ARP poisoning based MitM attack is analysed in [25].

A good overview of possible attacks can be found in [26] against the target protocol. The attacks are grouped into reconnaissance, attacks causing operational failures and attacks causing denial of service states. The effect of the different attacks is validated in a Hardware-In-the-Loop (HIL) Digital Station environment.

In the previous attack scenario, the authors assumed that the attacker already penetrated the SCADA system. [19] describes how this penetration can be done and used to inject IEC 104 packets into the system. The paper describes the necessary steps for a successful attack, like penetration, discovery, and injection.

From our point of view, the packet injection is an exciting step, but in that paper, they successfully injected commands to an IEC 104 server in a new session and not in an established session. To achieve it, they used a custom made tool to reset the session of the legitimate party (probably an operator) and initiated a new connection. Through the freshly built connection, they could issue IEC 104 commands to the server and control the grid that way. This kind of attack can be easily detected contrary to our proposed attack described in Section V.

A very detailed description of MitM attacks is presented in [20]. The authors of the paper tried two different attacks that required MitM position: replay attack and packet modification attack. In the replay attack, the attacker captured a packet containing an IEC 104 C_CS_NA_1 and replayed it without modification. This attack did not work as the sequence numbers in the TCP header were incorrect, resulting in packet drop and an alert in the Snort anomaly detection system. The second attack was more successful than the first one. This attack starts by waiting for a good message, capturing and modifying it (dropped the original packet), and forwarding the modified version. With this attack, they could control the power grid or cause false alarms for the operators. However, they needed valid packets to alter them in their own needs.

The authors of [21], [22], [23] did not focus on the IEC 104 protocol. Instead, they investigated what the possible effects of malicious NTP messages on a system using IEC 104 are. Some of the attacks may even affect IEC-62351 compliant systems (which proposes security improvements for the standard IEC 104) because it relies on NTP as well.

In [21] the authors reviewed if malicious time settings can be propagated across the entire network, and the results showed that IEC clients and servers started to use the malicious time in their ASDU messages. This behavior can be abused to de-synchronize control loops, and it even affects logging mechanisms in the system.

[22] focuses on causing de-synchronize just by manipulating the packet process rate of the NTP synchronization server's queue. This attack may be used to undermine the QoS of the communication.

In [23] the authors showed that in some cases, when the clock of communicating parties is not synchronized, it can even result in DoS attacks.

There are newer protocols to support the communication needs of power grids. One is the IEC 61850 [27] with security extensions defined in IEC 62351 [28]. However, these extensions are not flawless either. The authors in [29] showed three weaknesses in the protocol. The security mechanism of IEC 62351 can also be applied to the protocol in the focus of this paper, the IEC 104. This was demonstrated in a laboratory setup in [30] as well.

The security of power systems is a well-studied field, two great surveys are [31], [32]. A recent paper focusing on the security of substations using IEC 104, among other protocols, can be found in [33]. Readers interested in a greater picture of cyber-physical systems security are encouraged to check the survey in [8].

We used techniques from the previously mentioned papers and improved them to achieve better and stealthier attacks during our research. We first executed some of the attacks described in the previous papers. Preliminary results of this experiment were published in [34] by us. Then we moved to designing new attacks. We designed a new way that can be used for DoS attacks as well. Furthermore, we developed the packet injection attacks by injecting packets into an established session without resetting the original connection and without relying on valid packets. This method is far harder to detect because the communicating parties generate no log messages, and no invalid messages are sent, which could trigger the alert of an IDS system.

Before introducing our attack toolkit, we describe the potential entry points of a system in the next section.

## IV. ENTRY POINTS OF ATTACKS

In this section we highlight some key concepts of attack surfaces complementing our main topic. The basics and current state regarding operation of power system control and IT solutions can be found in literature (e.g. [35], [36], [37] )

Reliability is of main concern in every aspect of power system operation. Fault of a single element must be withstood without any disturbance (n-1 principle), be it a transmission line or a server. There is always a fallback procedure: the other transmission lines can take the extra load, a spare server jumps in. This mindset challenges the implementation of cyber-security measures. For example, an adaptive firewall cannot block traffic that is yet unknown - or else it can lead to emergency situations. This can also be exploited: a spare system or communication link can be quietly targeted and, at the right moment, locking out the primary system creates an automatic fallback to the infected system. Reliability also means that one cannot simply shut down or restart a critical system element, making it a harder and slower process to recover from an incident. The recovery process may also require significant human resource and time, especially if a black-out situation was created. This means that restoration of power supply can take several hours if no permanent damage or change was introduced to the system. If physical damage was done (e.g. circuit breakers), it can take weeks until the system could operate again in normal - n-1 withstanding - state.

Operation of power systems span across huge geographic areas involving multiple companies and personnel in different domains. The larger the system the larger the attack surface is, that malicious attackers might leverage. Generally there are several regulations that require making certain data publicly available which makes some areas of open source intelligence unavoidable (tenders, network development plans, transparency reports, press releases, etc.). Not only this, but the potential damage a successful attack can inflict makes power systems an attractive target. Figure 4 shows the typical areas where IEC 104 are used in system operations - focusing on SCADA systems.

Some example of potential weaknesses:

- Utilities must have periodic data exchange (schedules, measurements, grid models)
  - Corporate internal firewall - the interface must acquire data from operational databases inside SCADA systems
  - Data exchange process - corrupting or injecting false data
  - Using the sheer knowledge of the exchanged data can be leveraged at other attacks.
- Utilities have remote access to substations and sites
  - Suppliers' system - a vendor can have direct access to equipment and or software
  - Physically infiltrating a substation and attacking the control center from there - it can be easier than it sounds, considering a distant, unmanned substation, away from any settlement.
- Utilities rely on third parties
  - Equipment manufacturers or software developers can be leveraged - supply chain attack
  - Service providers - it can be e.g. ICT infrastructure as a service

From the many possibilities, some general key entry points can be pinpointed:

- The SCADA (Supervisory Control and Data Acquisition) system itself, which could potentially provide a complete control over the system. Due to this, it is one of the most protected system regarding perimeter defense. At the same time, it provides a relatively big attack surface due to required communication over a large geographic area.
- Several interfaces for the SCADA system are required to, for example, fulfill data exchange responsibilities and to keep the database up-to-date with a geographic information system (GIS). The attack surface here is twofold. On one hand, an attacker can inject false data into the external system, which could even have an impact that prevents the SCADA system operating. On the other hand, an attacker can exploit vulnerabilities of the several interfaces.
- Utility business process applications plays an important role in efficiency and organized operations. For example, a schedule from the power system market or a planned maintenance serve as a baseline of the network's planned state.

- Network equipment (routers, firewalls) is the backbone of most of the systems utilities have. At the same time, network equipment is typically used everywhere, not just in power systems. They tend to be cheaper and more accessible for hackers and security researchers. Therefore zero-day vulnerabilities are discovered in these much more often than in other pieces of equipment.
- RTU (Remote Terminal Unit) is the data concentrator for several equipment and provides protocol translation, routing and switching functionalities. These are found in every substation, and generally the primary and spare systems come from a different vendors. They even have primary and backup communication link that are physically independent from each other.
- Substation equipment connects to RTU and thus provides another way into the system. The attack surface is larger than that of RTUs, because there are usually two RTUs in a substation, while there are tens, hundreds of other equipment from several manufacturers.

The amount of equipment, cooperating parties and geographic distance coupled with the interdependence of physical and cyber space make the systematic listing of entry points and a complete risk assessment especially challenging. The general frameworks could be partly used with the electric power system like ISO 27000 series standard (e.g., 27019 for ICS in energy sector) or IEC 62443. In North America NIST also creates standards and guidelines - among others - for risk assessment (e.g. NIST SP 800-30) and for security of industrial control systems (e.g. NIST 800-82). A holistic and generalized framework that fits every purpose does not exist, however, the concepts and definitions could be extended and used to secure critical infrastructure systems. [38] shows an example for creating a combined approach for cyber-physical systems.

## V. Use cases for attacks against IEC 104

This section will introduce the attack vectors, the test environment, and the real-world devices used for evaluating the different attacks. We designed multiple attack scenarios based on our understanding of the protocol. Furthermore, we also tried attacks described in other papers, which were introduced in Section III.

The attacks were tested in a small simulator consisting of virtual machines and the OpenMUC [39] implementation of the IEC 104 protocol. The topology can differ for the scenarios, but it always had one IEC 104 server, one client, and at least one attacker.

The attacks are available for everyone wanting to reproduce our results. The attack scripts are publicly available on GitHub[2]. To make reproducing easier, we used Docker containers, and each attack scenario has a detailed step-by-step description. The topology of the attack is depicted on Figure 5.

---

[2]https://github.com/CrySyS/IEC-104-Attacks

Fig. 4. Generalized communication scheme of the power system control, highlighting usage of IEC 60870-5-104.



Fig. 5. The topology of the testbed.

As soon as we successfully attacked the purely software-based simulator, we also tested it on real devices, which included:

- an Opal-RT OP5707 based real-time simulator system. It can conduct hardware-in-the-loop simulation, which plays an important role in power system research, including communication-related studies.
- an Infoware RTU (MAB V2[3]). It is a piece of industrial-grade equipment produced by a local firm and is used widely in the power system of Hungary.
- a CitectSCADA based HMI by Schneider (Aveva). It is also a popular choice for substations' local HMI needs.
- a newer and an older version of EuroProt+ protection equipment by Protecta.

The most important attacks against the protocol are the following.

*A. Unauthorized access*

The protocol does not support authentication by design. Therefore if an attacker can connect to an IEC 104 server, they can execute arbitrary IEC 104 command on the server.

---

[3]http://www.infoware-zrt.hu/intelligens-rendszerek-mabv2-es-mab3-gwy-iranyitastechnikai-rendszer

We tried this attack against the OpenMUC implementation of the protocol and on two real devices (Infoware RTU, OpalRT simulator). As expected, none of them implemented extra authentication alongside the protocol. In the following, we will not discuss this attack because it is already explained in detail in [18].

### B. ASDU address field starvation

The idea of the attack is to use up every possible ASDU address field of the IEC 104 network. Since the ASDU address field is stored on 2 bytes, 65535 (actually 65533 since 0 and 65535 are not used) different addresses can be used. To deplete the address pool, the attacker must initiate connections to the IEC 104 server without closing the previously established connection. In theory, after the address pool becomes depleted, the server will not accept new connections.

The attack was carried out in the testbed with two attackers, one server, and one client. The OpenMUC implementation of the protocol can handle only 100 connections simultaneously. After reaching the limit, no new connections could be made to the server.

The attack was successful against real devices as well. The OpalRT simulator and the Infoware RTU could handle only a single connection. Therefore if an attacker could connect to the device, the operators could not interact with the server. This attack was also tested against two Protecta devices. Both could handle nine parallel connections. When all of them are acquired by an attacker, the operator would have no means of controlling the device. The small number of parallel connections can be easily exploited to prevent the operator from connecting to the server.

### C. TCP stream poisoning

The protocol uses TCP as the transport protocol. Therefore attacks targeting the TCP transport protocol can be used against the IEC protocol as well. An attacker who can inject packets in an active communication session can inject TCP packets with invalid checksums in the connection or inject a FIN or RST packet while impersonating a legitimate party of the communication to cause a denial of service. This attack was proven successful both in the simulator and in the case of real devices (between the Citect SCADA and the Infoware RTU, between the Opal-RT OP5707 and a RaspberryPi with OpenMUC) as well. TCP-based attacks will not be mentioned in the following because they are described in more detail in [18] and [19].

### D. Modifying IEC 104 sequence numbers

Each participant of the communicating parties keeps a record of the number of received and sent messages. Upon a message is received, these sequence numbers are validated. In the case of correct sequence numbers, the message is processed. Otherwise, the message is dropped, and the connection is terminated due to the incorrect sequence numbers. Since the protocol does not define a mechanism to protect the integrity of the message header, this behavior can be abused by an attacker who can modify the messages of active connection to cause a denial of service.

To modify the IEC 104 sequence numbers, the attacker needs to be in a MitM position to inspect and modify forwarded packets. We used *iptables* and the *NetfilterQueue* [40] package of *Python* to move traffic forwarding from kernel-space to user-space. Since, at the time of writing this paper, the *set_payload* function of NetfilterQueue was not working, we dropped the original packet and sent out a forged one using *Scapy* [41].

This attack was carried out with success in the simulator and with real devices (between the Citect SCADA and the Infoware RTU, between the Opal-RT OP5707 and a RaspberryPi with OpenMUC) as well. To achieve MitM position, we used *bettercap* [42], a framework for carrying out ARP poisoning attacks.

### E. Packet injection

The protocol does not define a mechanism for protecting the header or the body of the IEC messages. An attacker with man-in-the-middle capabilities can manipulate arbitrary packets of the communication. The attacker can even inject new packets into the communication. However, the attacker needs to be careful because the sequence numbers of the legitimate communicating parties will not match. To avoid the termination of the connection, the attacker needs to modify the sequence numbers to match the expectations of the destinations.

The setup used to inspect, modify or drop packets was the same as in V-D. Injecting a new packet to the communication and modification (in case of packet length changes) requires extra care because the TCP sequence numbers and the IEC 104 sequence numbers of the communicating parties will not match. Therefore to avoid the termination of the TCP session, it needs to be patched to match the receiver's expectations for all future packets as well. Therefore the attacker needs to keep a record of both session numbers on all of the attacked streams. The authors started to work on this topic and showed some preliminary results in [34]. This can results in a complete takeover of the communication, since the attacker can inject, modify or drop arbitrary packets without the possibility of trivial detection of the attack.

This attack was successful in the case of the simulator. Furthermore, it worked well in the case of real industrial devices (between the Citect SCADA and the Infoware RTU, between the Opal-RT OP5707 and a RaspberryPi with OpenMUC). To achieve MitM position, we used *bettercap* as in the previous attack.

The only remainder of an efficient attack is the targeting of the attack. The attacker knows what they want to achieve but does not know which ASDU address or IOA should be modified to achieve that goal. In the next section we show a method which can help the attacker in this question.

## VI. TARGETING ATTACKS IN UNKNOWN ENVIRONMENTS

### A. Introduction to the pairing based attack

Previously, we have shown numerous attacks against the IEC 104 protocol. These attacks included MitM attacks where

| $IUID$ | $T_1$ | $T_2$ | $T_3$ |
|--------|-------|-------|-------|
| 1000   | 8.7   | 11.2  | 9.3   |
| 2000   | 6.2   | 5.1   | 3.9   |

| $SID$ | $T_1$ | $T_2$ | $T_3$ |
|-------|-------|-------|-------|
| $A$   | 5     | 5     | 4     |
| $B$   | 8     | 12    | 10    |

the attacker could manipulate data transferred among the parties using IEC 104. However, the attacker could not decide what to manipulate because the pairing among real world identifiers (station IDs) and IEC 104 identifiers (ASDU Address; Information Object Addresses, IOA) is unknown. Therefore the attacker only could randomly set the values of different IOAs on different devices which could lead to harmful events, but not to full control.

In this part of the paper, we show how the attacker can learn the station ID-IOA pairs. The very same method can be used if the ASDU address is also unknown. In the following we will use the following nomenclature for the sake of simplicity:

SID:  station ID, real world identifier of a data point known to the attacker

IUID: IEC 104 unique identifier unknown to the attacker, ASDU address and IOA pair

The revealed pairing can help the attacker to successfully compromise and take over the control of the system. The attacker can use the previously described MitM position to eavesdrop and learn the values set by the IEC 104 client. This way, one can learn the values assigned to each IUID (remember, the ASDU contains the ASDU address, the IOA and the value itself). The attacker can also periodically collect some auxiliary information (e.g. query some public web page, or utilize other acquired intelligence sources) to learn the values assigned to each SID. The different sources of auxiliary information is described in Section VI-B. Based on the auxiliary information and the information collected at the MitM position the attacker can build an error matrix to find the most likely pairing. This algorithm is presented in the following artificially small example.

Table I. represent the data learned through eavesdropping the IEC 104 communication at the MitM position and collecting the auxiliary sources at times $T_1, T_2$ and $T_3$. For the sake of simplicity, only two IDs are used in the example. In real life scenarios, the number of IDs are much higher; the robustness of the algorithm with higher station numbers is analyzed in Section VI-F.

The next step of our algorithm is to calculate the error of the different pairings which is the sum of squared differences for every possible pairing. The possible pairings are the following ($SID - IUID$): A-1000, A-2000, B-1000, B-2000.

Based on the time series presented earlier, the calculations for the first pair would be the following:

$$ERROR(A, 1000) =$$
$$(8.7 - 5)^2 + (11.2 - 5)^2 + (9.3 - 4)^2 = 80.22$$

We can simply calculate the same error for every possible pair. The matrix of the errors can be represented as:

|      | $A$   | $B$   |
|------|-------|-------|
| 1000 | 80.22 | 1.62  |
| 2000 | 1.46  | 88.06 |

The last step of the algorithm is to choose pairs which covers the matrix and calculate the total error. It can be done in two ways in this example:

- 80.22 + 88.06 (SID=A is paired to IUID=1000 and SID=B is paired to IUID=2000)
- 1.46 + 1.62 (SID=A is paired to IUID=2000 and SID=B is paired to IUID=1000)

Our goal is to find the lowest one, so the second one will be the correct. From this, we can learn that the most likely pairing in this case was A-2000 and B-1000. Generally finding the best pairing is not so obvious. Fortunately the well-known Hungarian algorithm [43] solves exactly this problem.

The steps of the algorithm are the following:

1) Collect $n$ time series of IEC 104 communication for $n$ different IUIDs
2) Collect $n$ time series from auxiliary source for $n$ different SIDs
3) Calculate an $n$-by-$n$ matrix of squared differences of the time series
4) Find a minimal pairing based on the Hungarian algorithm
5) The result of the algorithm is the most likely IUID - SID pairing

We paired $n$ IUID to $n$ SID in the above example for the sake of simplicity. The very same algorithm can be used to pair $n$ by $m$ values ($n >= m$) if we have $m$ SIDs to pair with $n$ IUIDs.

### B. Auxiliary information sources

The pairing based attack requires an information source where the real station IDs (name of the variable) and related values can be found. Such an information source can be publicly accessible in many countries. These web pages display the current flows of different countries, stations or substations in near real time. Some examples are the following:

- Border crossing data per country and region: https://www.electricitymap.org
- EU detailed cross border physical flows and market data: https://transparency.entsoe.eu
- Real time electricity consumption and inter-regional flows from France: https://www.rte-france.com/en/eco2mix
- Hungarian Power System Actual Data: https://mavir.hu/web/mavir-en/hungarian-power-system-actual-data

The attacker can also use other private but not well protected information sources, like data series shared between different companies.

The finer the auxiliary information, the easier and more accurate the result is. In Section VI-F we show how the

accuracy of this information influences the results of the algorithm.

Transmission system (along with the transmission system operator - TSO) is most affected by this openness. The basic topology of the transmission system is relatively easy to reconstruct via satellite images, while for some regions it is even well known from publications (like the ENTSO-E materials in Europe). Generally, neighboring TSOs have a few connections with each other, thus - for example - the flows on the physical lines could be guessed from the data on energy exchange between countries. The transmission system has much less substations and transmission lines compared to the distribution system, which means that even little information defines the state of the system that adequately supplements the pairing algorithm.

*C. Implementation of the attack*

We have implemented the attack with 3 main parts, the IEC 104 eavesdropper, the auxiliary information collector (web crawler) and the SID-IUID match maker.

*1) IEC 104 eavesdropper:* The implementation of the eavesdropper is described in Section V. However in this particular case we only want to store the values of the IUIDs in time-series (time-IUID-value tuples). No modification or injection is needed in this phase.

*2) The web crawler:* The web crawler is similar to the IEC 104 eavesdropper as it downloads and stores the auxiliary information source in time-series (time-SID-value tuples).

*3) The match maker:* The match maker reads the stored data from the two sources and calculates the squared differences matrix. This matrix is the input of the Hungarian algorithm which calculates the pairs with the lowest error. The result is a list containing the SID-IUID pairs. The algorithm was implemented in Python using scipy and numpy.

*D. Verification of the matching algorithm*

The verification of the algorithm was done by testing it with many different cases. We could not use real life data in this research, therefore we wrote a Python script that generated numerous data series and also used a series from a simulator. We ran the algorithm on these artificial time-series and verified the results.

*E. Generation of data series*

The goal of this step is to generate data series that can be used for testing. We wanted to make the generated data series as realistic as possible, therefore the data was generated according to the following rules.

- The values of the stations are changing smoothly in real-life, therefore using plain random numbers would result in unrealistic scenarios. To avoid this, only the first value is random and the rest is generated using the ones before them.
- The publicly available web page nor the measurements themselves wouldn't represent accurate values in a real-life scenario, therefore some noise needs to be added to the data.

- The number of the values generated for each scenario is also random, therefore we can validate if it works properly with even a small amount of data available.

*F. Test scenarios*

We used a manually crafted test scenario (#5) and several artificially created base test cases (#1-#4) for performance evaluation. For the latter, the SID-IUID pairs are generated randomly for each test case. The results of the algorithm was compared to the ground truth. We ran millions of test cases with the following configurations.

Explanation of these constants are the following:

- CHANGE_RATE_MIN/CHANGE_RATE_MAX: A random number is generated in this interval, and the previous value is multiplied by it. This ensures that the values can only change smoothly.
- VALUE_COUNT_MIN/VALUE_COUNT_MAX: The number of values for each station is generated in this interval (length of test case).
- MIN_VALUE/MAX_VALUE: The value of the station is generated in this interval.
- ITERATION_COUNT: The number of test cases in the scenario.
- NOISE_MIN/NOISE_MAX: As mentioned earlier the website doesn't show accurate values, therefore the original value is multiplied by a number which is generated in the range of these numbers. This can be considered as quantization error.
- STATION_COUNT: The number of data points to pair.

The parameters of the five different scenarios are summarized on Table II.

TABLE II
PARAMETERS AND RESULTS OF THE SCENARIOS

| Parameter | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| CHANGE_RATE_MIN | | 0.8 | | | n.a |
| CHANGE_RATE_MAX | | 1.2 | | | n.a |
| VALUE_COUNT_MIN | 10 | 10 | 10 | 20 | 20 |
| VALUE_COUNT_MAX | 20 | 50 | 50 | 20 | 192 |
| MIN_VALUE | | -1000 | | | n.a |
| MAX_VALUE | | 1000 | | | n.a |
| ITERATION_COUNT | | $10^6$ | | | |
| NOISE_MIN | 0.7 | 0.8 | 0 | 0.9 | 0.85 |
| NOISE_MAX | 1.3 | 1.2 | 2 | 1.1 | 1.15 |
| STATION_COUNT | 13 | 13 | 13 | 100 | 13 |
| Success rate | 97.4% | 99.9% | 68.9% | 100% | 94.2% |

*1) Scenario 1:* In this scenario, the size of the data series was between 10 and 20, which is considered a rather small sample size. The algorithm correctly determined the pairs in 974827 cases out of 1000000. This means that in this case, the algorithm had a success rate of 97.4%.

*2) Scenario 2:* In this scenario, we tested the effect of longer time series (10 to 50 values per station) with smaller noise. We expected to get better results than in the previous scenarios as more input data was used in the pairing phase. The algorithm correctly determined the pairs in 999282 cases out of 1000000. This means that a success rate of 99.9% was achieved.

*3) Scenario 3:* In this scenario, the noise rate was enormously large (in a real scenario, it would be less than 10%, here we have 100%), and the value count per station was between 10 and 50. The algorithm still determined the pairs correctly in 689448 cases out of 1000000. This means that in this case, the algorithm had a success rate of 68.9%.

*4) Scenario 4:* In this scenario, the algorithm correctly determined the pairs in 1000000 cases out of 1000000. This means that in this case, the algorithm had a success rate of 100%. In this scenario, the number of stations was greatly increased, but the run-time of the matching scaled linearly. This indicates that the algorithm is suitable for larger sets with reasonable running times. We think this scenario is the most realistic in terms of parameter selection.

*5) Scenario 5:* This scenario significantly differs from the previous ones. The values for this scenario was crafted manually to mimic real-life values. Creating a lifelike simulated environment is generally challenging, however, here the goal was to simply have a system which has the same correlation behavior as a real one: the different measurements on substations (i.e. a node in a graph) and on transmission lines (i.e. an edge in a graph) all have a dependency on each other. The strongest connection among these measurements comes from the physical properties of a network which is represented by load-flow calculation that is generally used for these purposes in the industry. A correlation - although to a much lesser extent - also could exist via the natural shape of load profiles; this approach is also widely used in the industry since decades. For this specific scenario, consumption and generation values were based on actual measurements and statistically created profiles. Moreover, a simulated network topology was used to estimate the current and voltage values for every time step with load-flow calculations. Thus, even the effect of one node on an other is also apparent. Also, a variety of load behaviour was represented (photovoltaic, wind, gas power plant; industrial, commercial and aggregated residential load). We created a single set of timeseries as a foundation and added random noise (to represent measurement uncertainty) just like with the other scenarios and used a variety of time periods to represent that an attacker starts to eavesdrop at a random time. The algorithm correctly determined the pairs in 942000 cases out of 1000000. This means that in this case, the algorithm had a success rate of 94.2%. It is important to note that there are periods when some values remain the same for some nodes, which makes pairing really hard for the algorithm (e.g. photovoltaic plants are typically not producing anything for several hours during the night resulting in long, indistinguishable sequences of zeros).

*G. Evaluation*

The previous test scenarios showed that the algorithm works well even with extreme configurations. Scenario 4 was closest to reality in terms of complexity. In that case, the algorithm had a success rate of 100%. The data values were as realistic as possible in Scenario 5, and that scenario also succeeded over 94% of the cases. Therefore we affirm that this approach can be used to match SIDs to IUIDs. The speed of the algorithm is

also acceptable, Scenario 4 had the longest run-time, but it still took less than 1 second to calculate one case. We argue that the described algorithm can be successfully used as an early phase of an attack, where the target IUIDs can be discovered. The algorithm should run only once after a long data collection phase so the running time of the algorithm is not so important.

As the results show, the pairing works very well, but a successful attack chain requires many steps described in this paper. First, the attacker needs to access the network where the IEC communication is used (the possible entry points are analysed in Section IV). Then the attacker needs to get in a position where they can eavesdrop on the IEC server and client (in this step, the attacker gathers identifier and value pairs, the possible methods are introduced in Section V). Meanwhile, eavesdropping on the transmission, the attacker also needs to find a reliable data source where the station's values and names are present. After the attacker has gathered enough information, they can run the Hungarian algorithm to find the correct SID and IUID pairs as described in this section. When the pairing is done, the attacker can start various attacks against the power grid to control it as they desire. With this knowledge, masquerading and precisely targeting the attack becomes feasible. For example, the adversary could open circuit breakers without the system operators knowing it, thus leaving specific areas without power. Or, the adversary could remain hidden until it modifies the network enough to cause a large scale blackout.

The risk of the whole scenario can be analysed based on the DREAD scoring system [44], where the components of the risk are Damage potential, Reproducibility, Exploitability, Affected users, and Detectability. The scoring systems uses High, Medium and Low scores. The Damage potential of the scenario is high as physical damage even causing fire at the substation is probable. The Reproducibility of the attack is medium as special knowledge of power systems and network protocols are also required at the same time. The Exploitability is Medium as a vulnerable entry point must be found at the beginning, but after that the attack is straightforward. The number of Affected Users can be very high as a blackout around a substation can affect thousands of people. The attack can be detected relatively easily as it creates different anomalies in the operation. A well-configured anomaly detector can detect it and trigger an alarm. This means that the risk coming from the detectability is low. A successful attack scenario can violate all components of the CIA objectives: the Confidentiality of identifiers used inside the substation is violated by the pairing. The Integrity of the commands and measured values is affected by the selective modification of the IEC 104 messages, while the Availability of the power system might be violated by selectively opening and closing circuit breakers.

If someone wants to avoid such attacks, a systematic approach must be used. It starts with the application of the risk assessment frameworks discussed in Section IV and the realization of countermeasures mentioned in the frameworks and in Section III.

## VII. Summary

In this paper, we analyzed the security of the IEC 104 protocol. We showed existing attacks against the protocol and also designed new attack vectors. To supplement, we published our attack scripts to make it possible for everyone to reproduce our attacks.

We assembled a realistic attack scenario in a holistic approach comprising the following parts:

- We enumerated the possible entry points a threat actor can use to access IEC 104 communication.
- We showed and created new techniques which can result in man in the middle position for the attacker and showcased the operation with actual equipment.
- By collecting traffic from the attacked network and from auxiliary sources we demonstrated how the real identity of the information object addresses can be revealed.

By leveraging the demonstrated methods, one can get into a position from where taking over the entire telecontrol system is straightforward.

In the future, we will focus on extending the framework with advanced false packet creation methodologies which can increase the chance that the attack itself could remain hidden. This involves synchronised packet injection at multiple places that produces a consistent state of the network. It would also be interesting to test our attack scenario in a real deployment (especially the accuracy of the pairing), but it is challenging to persuade a TSO to let us access their internal systems.

## References

[1] "IEC 60870-5-104, part 5-104: Transmission protocols – network access for iec 60870-5-101 using standard transport profiles," 2006.

[2] "IEC 62351, power systems management and associated information exchange - data and communications security," 2021.

[3] R. M. Lee, M. J. Assante, and T. Conway, "Analysis of the cyber attack on the ukrainian power grid," *Electricity Information Sharing and Analysis Center (E-ISAC)*, vol. 388, 2016.

[4] D. E. Whitehead, K. Owens, D. Gammel, and J. Smith, "Ukraine cyber-induced power outage: Analysis and practical mitigation strategies," in *2017 70th Annual Conference for Protective Relay Engineers (CPRE)*. IEEE, 2017, pp. 1–8. [Online]. Available: **DOI**: 10.1109/CPRE.2017.8090056

[5] L. Constantin, "Solarwinds attack explained: And why it was so hard to detect," 2020. [Online]. Available: https://www.csoonline.com/article/3601508/solarwinds-supply-chain-attack-explained-why-organizations-were-not-prepared.html

[6] F. Wortley, C. Thompson, and F. Allison, "Log4Shell: RCE 0-day exploit found in log4j 2, a popular Java logging package," 2021. [Online]. Available: https://www.lunasec.io/docs/blog/log4j-zero-day/

[7] E. Burke, "Targeted cyberattack takes out Vodafone Portugals," 2022.

[8] I. Zografopoulos, J. Ospina, X. Liu, and C. Konstantinou, "Cyber-physical energy systems security: Threat modeling, risk assessment, resources, metrics, and case studies," *IEEE Access*, vol. 9, pp. 29 775–29 818, 2021. [Online]. Available: **DOI**: 10.1109/ACCESS.2021.3058403

[9] P. Matoušek, "Description and analysis of iec 104 protocol," *Faculty of Information Technology, Brno University of Technology, Tech. Rep, 2017*. [Online]. Available: https://www.fit.vut.cz/research/publication/11570

[10] E. Hodo, S. Grebeniuk, H. Ruotsalainen, and P. Tavolato, "Anomaly detection for simulated IEC-60870-5-104 traffic," in *Proceedings of the 12th International Conference on Availability, Reliability and Security*, 2017, pp. 1–7. [Online]. Available: **DOI**: 10.1145/3098954.3103166

[11] C.-Y. Lin and S. Nadjm-Tehrani, "Timing patterns and correlations in spontaneous SCADA traffic for anomaly detection," in *22nd International Symposium on Research in Attacks, Intrusions and Defenses ({RAID} 2019)*, 2019, pp. 73–88. [Online]. Available: https://www.usenix.org/conference/raid2019/presentation/lin

[12] Y. Yang, K. McLaughlin, T. Littler, S. Sezer, B. Pranggono, and H. Wang, "Intrusion detection system for IEC 60870-5-104 based SCADA networks," in *2013 IEEE power & energy society general meeting*. IEEE, 2013, pp. 1–5. [Online]. Available: **DOI**: 10.1109/PESMG.2013.6672100

[13] M. A. S. Arifin, D. Stiawan, J. Rejito, M. Y. Idris, R. Budiarto et al., "Denial of service attacks detection on scada network iec 60870-5-104 using machine learning," in *2021 8th International Conference on Electrical Engineering, Computer Science and Informatics (EECSI)*. IEEE, 2021, pp. 228–232. [Online]. Available: **DOI**: 10.23919/EECSI53397.2021.9624255

[14] M. Anwar, A. Borg, and L. Lundberg, "A comparison of unsupervised learning algorithms for intrusion detection in iec 104 scada protocol," in *2021 International Conference on Machine Learning and Cybernetics (ICMLC)*. IEEE, 2021, pp. 1–8. [Online]. Available: **DOI**: 10.1109/ICMLC54886.2021.9737267

[15] S. V. B. Rakas, M. D. Stojanović, and J. D. Marković-Petrović, "A review of research work on network-based scada intrusion detection systems," *IEEE Access*, vol. 8, pp. 93 083–93 108, 2020. [Online]. Available: **DOI**: 10.1109/ACCESS.2020.2994961

[16] O. A. Alimi, K. Ouahada, A. M. Abu-Mahfouz, S. Rimer, and K. O. A. Alimi, "A review of research works on supervised learning algorithms for scada intrusion detection and classification," Sustainability, vol. 13, no. 17, p. 9597, 2021. [Online]. Available: **DOI**: 10.3390/su13179597

[17] P. Matoušek, O. Ryšavỳ, and M. Grégr, "Increasing visibility of IEC 104 communication in the smart grid," in *6th International Symposium for ICS & SCADA Cyber Security Research 2019 6*, 2019, pp. 21–30. [Online]. Available: **DOI**: 10.14236/ewic/icscsr19.3

[18] P. Radoglou-Grammatikis, P. Sarigiannidis, I. Giannoulakis, E. Kafetzakis, and E. Panaousis, "Attacking IEC-60870-5-104 SCADA systems," in *2019 IEEE World Congress on Services (SERVICES)*, vol. 2642. IEEE, 2019, pp. 41–46. [Online]. Available: **DOI**: 10.1109/SERVICES.2019.00022

[19] Q. S. Qassim, N. Jamil, M. Daud, N. Ja'affar, S. Yussof, R. Ismail, and W. A. W. Kamarulzaman, "Simulating command injection attacks on IEC 60870-5-104 protocol in SCADA system," *International Journal of Engineering & Technology*, vol. 7, no. 2.14, pp. 153–159, 2018. [Online]. Available: **DOI**: 10.14419/ijet.v7i2.14.12816

[20] P. Maynard, K. McLaughlin, and B. Haberler, "Towards understanding man-in-the-middle attacks on IEC 60870-5-104 SCADA networks," in *2nd International Symposium for ICS & SCADA Cyber Security Research 2014 (ICS-CSR 2014) 2*, 2014, pp. 30–42. [Online] Available: **DOI**: 10.14236/ewic/ics-csr2014.5

[21] A. Baiocco and S. D. Wolthusen, "Causality re-ordering attacks on the IEC 60870-5-104 protocol," in *2018 IEEE Power & Energy Society General Meeting (PESGM)*. IEEE, 2018, pp. 1–5. [Online]. Available: **DOI**: 10.1109/PESGM.2018.8586010

[22] J. Wright and S. Wolthusen, "Time accuracy de-synchronisation attacks against IEC 60870-5-104 and IEC 61850 protocols," in *2019 IEEE Power & Energy Society Innovative Smart Grid Technologies Conference (ISGT)*. IEEE, 2019, pp. 1–5. [Online]. Available: **DOI**: 10.1109/ISGT.2019.8791558

[23] A. Baiocco and S. D. Wolthusen, "Indirect synchronisation vulnerabilities in the IEC 60870-5-104 standard," in *2018 IEEE PES Innovative Smart Grid Technologies Conference Europe (ISGT-Europe)*. IEEE, 2018, pp. 1–6. [Online]. Available: **DOI**: 10.1109/ISGTEurope.2018.8571604

[24] J. Jarmakiewicz, K. Parobczak, and K. Maślanka, "Cybersecurity protection for power grid control infrastructures," *International Journal of Critical Infrastructure Protection*, vol. 18, pp. 20–33, 2017. [Online]. Available: **DOI**: 10.1016/j.ijcip.2017.07.002

[25] D. Deb, S. R. Chakraborty, M. Lagineni, and K. Singh, "Security analysis of mitm attack on scada network," in *Machine Learning, Image Processing, Network Security and Data Sciences: Second International Conference*, MIND 2020, Silchar, India, July 30-31, 2020, Proceedings, Part II 2. Springer, 2020, pp. 501–512. [Online]. Available: **DOI**: 10.1007/978-981-15-6318-8_41

[26] L. Erdődi, P. Kaliyar, S. H. Houmb, A. Akbarzadeh, and A. J. Waltoft-Olsen, "Attacking power grid substations: An experiment demonstrating how to attack the scada protocol iec 60870-5-104," in *Proceedings of the 17th International Conference on Availability, Reliability and Security*, ser. ARES '22. New York, NY, USA: Association for Computing Machinery, 2022. [Online]. Available: **DOI**: 10.1145/3538969.3544475

[27] R. E. Mackiewicz, "Overview of IEC 61850 and benefits," in *2006 IEEE Power Engineering Society General Meeting*. IEEE, 2006, pp. 8–pp. [Online]. Available: **DOI**: 10.1109/PSCE.2006.296392

[28] S. S. Hussain, T. S. Ustun, and A. Kalam, "A review of IEC 62351 security mechanisms for IEC 61850 message exchanges," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 9, pp. 5643–5654, 2019. [Online]. Available: **DOI**: 10.1109/TII.2019.2956734

[29] M. Strobel, N. Wiedermann, and C. Eckert, "Novel weaknesses in IEC 62351 protected smart grid control systems," in *2016 IEEE International Conference on Smart Grid Communications (SmartGridComm)*. IEEE, 2016, pp. 266–270. [Online]. Available: **DOI**: 10.1109/SmartGridComm.2016.7778772

[30] M. G. Todeschini and G. Dondossola, "Securing IEC 60870-5-104 communications following IEC 62351 standard: lab tests and results," in *2020 AEIT International Annual Conference (AEIT)*, 2020, pp. 1–6. [Online]. Available: **DOI**: 10.23919/AEIT50178.2020.9241101

[31] W.-W. Li, W.-X. You, and X.-P. Wang, "Survey of cyber security research in power system," *Power System Protection and Control*, vol. 39, no. 10, pp. 140–147, 2011.

[32] W. Wang and Z. Lu, "Cyber security in the smart grid: Survey and challenges," *Computer networks*, vol. 57, no. 5, pp. 1344–1371, 2013. [Online]. Available: **DOI**: 10.1016/j.comnet.2012.12.017

[33] S. Hussain, J. H. Fernandez, A. K. Al-Ali, and A. Shikfa, "Vulnerabilities and countermeasures in electrical substations," *International Journal of Critical Infrastructure Protection*, p. 100 406, 2021. [Online]. Available: **DOI**: 10.1016/j.ijcip.2020.100406

[34] P. György and T. Holczer, "Attacking IEC 60870-5-104 protocol," in *Proceedings of the 1st Conference on Information Technology and Data Science*, 2020, pp. 140–150. [Online]. Available: http://ceur-ws.org/Vol-2874/paper13.pdf

[35] M. S. Thomas and J. D. McDonald, *Power System SCADA and Smart Grids*. CRC press, 2017. [Online]. Available: **DOI**: 10.1201/b18338

[36] K. Siozios, D. Soudris, D. Anagnostos, and E. Kosmatopoulos, *IoT for Smart Grids: Design Challenges and Paradigms (Power Systems)*. Springer, 2019. [Online]. Available: **DOI**: 10.1007/978-3-030-03640-9

[37] S. K. Khaitan, C.-C. Liu, and J. D. McCalley, *Cyber Physical Systems Approach to Smart Electric Power Grid*. Springer, 2015. [Online]. Available: **DOI**: 10.1007/978-3-662-45928-7

[38] S. C. Matta, A. M. Shaaban, C. Schmittner, A. Pinzenöhler, E. Szalai, and M. Tauber, "Risk management and standard compliance for cyber-physical systems of systems," *Infocommunications Journal*, vol. 13, no. 2, pp. 32–39, 2021. [Online]. Available: **DOI**: 10.36244/ICJ.2021.2.5

[39] "IEC 60870-5-104," https://www.openmuc.org/iec-60870-5-104/, 2019.

[40] "NetfilterQueue," https://pypi.org/project/NetfilterQueue/, 2017.

[41] "Scapy," https://scapy.net/, 2021.

[42] "bettercap," https://www.bettercap.org/, 2021.

[43] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955. [Online]. Available: **DOI**: 10.1002/nav.3800020109

[44] J. Meier, *Improving web application security: threats and counter-measures*. Microsoft press, 2003.

**János Csatár** received the Ph.D. degree in Electrical Engineering from the Budapest University of Technology and Economics (BME) in 2019. Since 2019 he has been working as an assistant professor in the Smart Power Laboratory (SPL), Department of Electrical Power Engineering, Budapest University of Technology and Economics. Fields of interest: His research interest for the PhD was power system modeling with a special focus on distribution system. Recently his research interest focuses on hardware-in-the-loop simulation, co-simulation and interdependence of power systems and communication networks.

**Péter György** holds a Master's degree from the Budapest University of Technology and Economics (BME), which he obtained in 2022. He has actively participated in Capture The Flag competitions since 2019 as a member of the c0r3dump team. In 2021, Péter joined Ukatemi Technologies Plc, where he currently serves as the leader of the penetration testing team. Fields of interest: Péter focused his research on identifying vulnerabilities through black-box fuzzing during his Master's program. However, his current interests primarily revolve around network penetration testing and exploring various facets of web security.

**Tamás Holczer** received the Ph.D. degree in Computer Science from the Budapest University of Technology and Economics (BME) in 2013. Since 2013 he has been working as an assistant professor in the Laboratory of Cryptography and System Security (CrySyS), Department of Telecommunications, Budapest University of Technology and Economics. Fields of interest: In the past his research interests and his Ph.D. dissertation were focused on the privacy problems of wireless sensor networks and ad hoc networks. Lately he is working on the security aspects of cyber physical systems. The research topics include: security of industrial control networks, honeypot technologies in embedded systems, network monitoring and intrusion detection in industrial networks, and security aspects of intra-vehicular networks.

# Evaluation of Bone Conduction and Active-Noise-Cancellation Headsets Based on Listening Tests in a Virtual Environment

György Wersényi

*Abstract*—Alternative design headsets incorporating active-noise-cancellation or bone conduction were evaluated in listening tests in a virtual reality environment. Virtual sound sources in the horizontal plane had to be identified using stereo panning in the frontal hemisphere. In addition, transfer characteristics and damping effects were measured with a dummy-head. Results indicate that up to five source locations can be used in real applications with high accuracy in virtual scenarios, independent of the spectral content of the excitation signals. Furthermore, the use of noise cancellation in presence of 80 dB background noise does not improve performance. Commercially available bone conduction headsets can provide the same detection accuracy even if the subjective sound quality is lower.

*Index Terms*—Active Noise Cancellation, headphone, measurement, bone conduction, virtual reality

## I. INTRODUCTION

Localization is the procedure during which human subjects try to find sound sources [1], [2]. Directional hearing is based on interaural time and intensity differences between the signals on the eardrums, as well as on spectral filtering cues of the outer ears [3]. Furthermore, many other parameters influence the localization accuracy, i.e., spectral and temporal properties of the signal, experimental environment (real-life or virtual), playback devices, additional sound signal processing, presence of distractor and masking sounds, and experience of the subjects, etc. It is also an important parameter whether source locations must be "pointed to" in a so-called absolute localization task, or subjects must identify source locations from a limited number of possibilities (e.g., a given number of directions or loudspeakers). Both methods can be applied in real and virtual environments, and the latter is usually a simplified task due to the limited possibilities. The term lateralization can be also used along with localization, if the effect of externalization plays a significant role [4].

The significance of virtual reality (VR) environments has increased with the introduction of modern playback and feedback devices for both desktop and for full immersive VR with headsets. Furthermore, virtual audio displays (VAD) play a significant role in various applications from assistive technologies to simulators and gaming, extending or replacing visual screens [5], [6]. It is often a requirement to have privacy and to restrict the virtual audio/visual experience to the user only. On the other hand, safety issues - especially in mobile applications - may require some contact with the environment regarding environmental sounds, such as traffic noise, alarm sounds, and speech communication. Traditional headphones covering the ears or plugged in the ear canal damp the outer world, although the hearing modality may be the only sense receiving information from the surrounding environment, if vision is focused elsewhere. VR environments and simulation tools exist not only to create 3D visual spaces, but they are also appropriate for audio scene rendering, usually over two-channel headsets. Various methods can assist the simulation of directional information of virtual sources. Virtual labs for scientific purposes offer solutions for individual needs for experiments [7]–[9].

The introduction of nontraditional headsets with different solutions to enhance sound quality, can also provide convenience services or extended safety also for VADs and VR scenarios. Bone conduction is a technique that uses the human skull to transmit acoustic vibrations while leaving the ear canals open [10], [11]. Another interesting and widespread method for high quality headphones is the presence of an active-noise-cancellation circuit to increase performance [12], [13]. Localization, however, may be influenced and is different in contrast to traditional headphones. The next subsections discuss these two techniques briefly.

### A. Bone Conduction

Bone conduction (BC) headsets or "bonephones" may be an alternative to traditional headphones and earphones if an open ear canal entrance is required during operation. The most important application areas include electronic travel aids (ETAs) and wearables for assistive technology (mostly for the visually impaired [14], [15]); outdoor sport activities (safety during running or biking in traffic); military or combat activities with sparsely occurring commands [16]; as well as any application where environmental sounds and noises can not be blocked [17]. Although VR applications usually try to block the outer world for full immersion, many applications can benefit from an open ear canal. In case of BC, sound quality is inferior to traditional headphones due to the indirect sound path through the skull instead of the direct sound path through the eardrum. A restricted frequency range is common, while directional information can be disturbed. Furthermore, transmission depends strongly on the transducer position and whether the ear canal is open or closed [18]–[20].

Széchenyi István University, Győr, Hungary (e-mail: wersenyi@sze.hu).

Higher sensitivity (lower hearing threshold), however, can occur during skin penetration [21]. A threshold in case of a plugged ear canal was found to be 10-20 dB lower than the unplugged threshold at low frequencies (200 Hz and below), beginning to converge with the unplugged threshold at medium frequencies (250 – 1000 Hz), and the two being equal above 2 kHz [22].

The mastoid would be a preferable transducer location relative to the forehead or temple because it contains the inner ear. It is relatively immune to the interference associated with muscle tissue operating the jaw, and it allows stereo presentation of sounds [23]. In a study, binaural hearing ability of normal hearing adults with BC and air conduction (earphones) was contrasted. BC was applied in the audiometric position on the mastoid [24]–[27]. Results confirmed that binaural hearing processing with bilateral BC stimulation is present. However, the binaural benefit was overall greater with air conduction stimulation. On the other hand, other studies reported no significant difference [16], [28], [29]. Nevertheless, current devices do not use the mastoid, but the jaw bone and the zygomatic (cheek) bones in front of the ear.

BC use is traditionally limited to monaural applications due to the high propagation speed of sound in the human skull. Spatial audio does not occur naturally through bone conduction, although interaural level and intensity differences can be simulated. It was shown that stereo bone conduction headsets can be used to provide a limited amount of interaural isolation in a dichotic speech perception task [30]. The results suggested that reliable spatial separation is possible with bone conduction headsets, but they probably cannot be used to lateralize signals to extreme left or right apparent locations. However, the degree of lateralization can be similar to that of produced by using headphones. Results from an empirical user study conducted to compare one BC device, headphones, and a speaker array showed that subjects performed the best by using physical speakers with stationary sounds. Nevertheless, there was no difference in accuracy between the speakers and the BC device (outperforming even standard headphones) for moving sounds [31]. Elevation cues can also be adjusted for 3D simulations using BC [32]. Monaural spectral cues are also responsible for front-back and up-down discrimination and externalization of sound sources, and problems occurring frequently during headphone playback. Elevation cues can be simulated for BC with different methods using Head-Related Transfer Functions (HRTF) or simple high-pass and low-pass filtering [15], [33], [34].

### B. Active Noise Cancellation

Active-Noise-Cancellation (ANC) is a technique, where an incident, unwanted sound signal is "cancelled out" by inverting and adding it to the computer-generated signal based on interference. In case of a perfect destructive interference, the sum of the incoming and inverted signals equals zero. Although the theory is very simple, practical realization of ANC circuits in noise cancelling headphones face many problems [35]–[37]. Some of the ANC headsets offer not only cancelling, but also allow the user to hear environmental noises by reducing isolation or enabling communication with the environment by pushing a button (and avoiding taking off the headphone), i.e., via pass-thru, mic-thru or hear-thru functions. This is beneficial for augmented reality applications as well, where environmental sounds and computer-generated signals are mixed together.

ANC, in general, is supposed to increase damping, allowing for reduced loudness during playback and leveling the subjective quality of the playback. This reduction of disturbing noises may also result in better localization performance in listening tests.

This paper presents the results of a classic subjective listening test, using BC and ANC headsets in a virtual environment where virtual sound sources had to be identified using various excitation signals. The goal was to test whether these cost-effective, commercially available devices can be an alternative to traditional headphones in applications where directional information is important. Section 2 presents the measurement setup, including the virtual environment, the headsets and the experimental procedure. The subsequent section presents the results of the listening tests and the objective measurement of the technical parameters (transfer function and damping). Then, the results will be discussed and conclusions will be drawn together by highlighting future research directions.

## II. Measurement setup

### A. The virtual environment

The virtual environment was created by applying the MaxWhere platform. This multipurpose 3D collaborative environment is a versatile platform made up of various virtual spaces that can be used effectively in education, virtual laboratory tests, and even for testing memory capabilities [38], [39].

Upon initialization, the user could set personal data (name/ID, age, gender) and the simulation environment. The main settings included the number of virtual sound source directions (3, 5 or 7), and the type of the sound sample (impulse, 1 kHz sinus, white noise, and female speech). Furthermore, a checkbox was marked if the user was visually impaired, if the ear canal entrance was closed, and in case of a looped signal presentation. Currently, only sighted subjects participated, thus the first checkbox was always set to false. The ear canal was closed in case of regular headphones and it was left open in case of BC devices. By default, all sound samples were played back once without being looped.

Figure 1 shows the directions of the virtual sources. Figure 2 depicts the screenshot of the arrangements where subjects had to click. In case of three directions, the front $0°$, the left $+90°$ and the right $-90°$ were set. Because left and right correspond to a radiating sound only from the left or the right speaker respectively, the actual vizualisation would require a more complicated way where the user had to turn. To avoid this (that would make the usability worse), the extreme left and right speakers were displayed relative to the listeners spot (marked red). The users did not have any problem to interpret the user interface.

Evaluation of Bone Conduction and Active-Noise-Cancellation
Headsets Based on Listening Tests in a Virtual Environment



Fig. 1. Arrangement of the virtual speakers in three, five, and seven directions.



Fig. 2. Screenshots of the virtual environments. After playback of the sound, subjects identify the sound source by clicking on the loudspeaker icon.

In the case of five locations, two additional locations of 45° "halfway" left and right from the frontal direction were included. In the case of seven directions, instead of 45°, 30° and 60° were introduced. Directions were simulated by applying a simple two-channel stereo intensity panning between the left and the right channels. Although two-channel virtual reproduction of sound scapes often suffer from in-the-head localization (lack of externalization), complex full binauralization is not necessarily required for applications where usability does not rely on externalization.

The excitation signals are the following:

- computer-generated white noise sample of 1.5 sec with a flat spectrum between 20 Hz and 20 kHz;
- a female speech sample with a base frequency of 430 Hz and harmonic components up to 8.2 kHz (2 sec);
- 1 kHz sine, with damping of the higher harmonics more than 60 dB (1.5 sec);
- an impulse-like complex sound generated from a 258 Hz base frequency signal and high order odd harmonics up to 15 kHz (2 sec).

During a test session, all sound sources were activated three times in a randomized order (unknown to the user). Thus, in case of three directions, the users had to identify a source location 9 times, in case of five sources it was 15, and in case of seven locations the experiment lasted for 21 clicks. After

|  | BC | | ANC | | | | | | | |
|  | AudioBone | AfterShokz | Bose | | | Sennheiser | | | |
|  | NOISE OFF | NOISE OFF | ANC OFF NOISE OFF | ANC ON NOISE ON | ANC OFF NOISE ON | ANC OFF NOISE OFF | ANC ON NOISE ON | ANC OFF NOISE ON |  |
|---|---|---|---|---|---|---|---|---|---|
| Impulse | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 24 |
| Speech sample | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 24 |
| 1 kHz sine | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 24 |
| White noise | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 3/5/7 | 24 |
| Nr. of test/session | 12 | 12 | 12 | 12 | 12 | 12 | 12 | 12 |  |
| Nr. of test/subject | 24 | | 72 | | | | | | |
| Nr. of subjects | 25 | | 25 | | | | | | 50 |
| Total nr. of tests | 600 | | 1800 | | | | | | 2400 |

Fig. 3. Overview of the experiments. Red numbers indicate the number of sound directions.

completing the task, a log file in the JSON format was saved with the personal data, time and date (timestamp), headphone type, completion time in seconds and the total points obtained. The maximum points were 9, 15 and 21 respectively, and the result was also calculated in percentage. Furthermore, the number of correct and incorrect identifications for each location was recorded.

*B. Experimental setup*

Laboratory sessions were conducted in a silent, acoustically controlled room. The room had originally been designed for sound recordings with insulated walls and a reverberation time below 1,5 s. A loudspeaker was also installed in front of the listener behind a desk where the computer was placed.

The subject was sitting in a chair and he was controlling the experiment by mouse clicks on a standard laptop screen. In the case of a distractor sound, the loudspeaker radiated a previously pre-recorded traffic noise, having a SPL level of 80 dB at the listening point. An external distractor could be used to test masking effects and/or the effect of active noise cancellation.

Prior to the first session, the measurement procedure and the goal of the experiment were described to the subjects, which was followed by a brief introduction to the sound signals and to the user interface (how to click, how to get feedback etc.). During the experiment, subjects wore one of the headphones. Headsets were placed on the head and adjusted to fit the most comfortably. The users set the loudness to the most comfortable level, neither too loud, nor too quiet. Due to the relatively low sensitivity of the BC types, the maximum volume was set on the laptop to provide appropriate loudness levels. The vibrational actuators were positioned on the jawbones in front of the ear, instead of on the mastoid bones, similarly to [31]. Finally, personal data were entered for the JSON files.

One session was to test one of the headsets on all three scenarios with four signal types. This resulted in 12 test runs in total. The test always started with the impulse (3, 5 and 7 directions), followed by speech, sinus, and white noise. The total time needed for one session was about 30 minutes. The procedure was repeated with the other headset. In the case of BC devices, 25 subjects evaluated them in two sessions. In the case of ANC devices, a different set of 25 participants evaluated them in six sessions. All 50 participants had normal hearing self reportedly (no audiometric screening was applied). In the absence of distractor noise, the measurement was conducted with ANC off only. Figure 3 shows an overview of the sessions.

### C. Devices

For the experiment, two types of BC and two types of ANC headsets were used as listed (Figure 4).

- Aftershokz Sportz M3 (wired) [34], [40];
- AudioBone Deluxe GDP 02 [31], [41];
- Sennheiser PXC 450 [42];
- Bose QC25 [43].

The AudioBone was already used in other experiments and the AfterShokz models are one of the most popular devices nowadays (usually for outdoor sport activities). Both are cost-efficient and they can be placed on the jawbone instead of behind the ears. ANC models are more widespread and commercially available from different vendors in various forms (supraaural, in-ear and buds).

Transfer characteristics of the ANC devices were measured following the standard protocol for headphone measurements [44], [45]. In a semi-anechoic room, the Brüel and Kjaer Type 5128 head and torso simulator was connected to the PULSE data acquisition system. The same headphone was placed and replaced 10 times for averaging, and two-channel measurements using white noise excitation, which were performed on



Fig. 4. Four headsets used in the experiments. AfterShokz (top left), AudioBone (top right), Bose QC25 (bottom left), and Sennheiser PXC 450 (bottom right).

the left and the right ears, respectively. Test were performed both with and without ANC [46].

### III. RESULTS

#### A. Technical parameters

In the case of BC devices, the conventional transfer function measurements designed for headphones could not be used. The output signal is not the sound pressure on the eardrums, but the vibrations on the skull and/or inside the head (in the middle and inner ear). The sensation is very complex based on various transmission paths of vibrations with an airborne sound through the ear canal being also present. There is no standard method to determine the transfer characteristics. Although there is a measurement equipment called artificial mastoid for the calibration of audiometric bone vibrators, it was designed in a very simplified form to model the human receiver and for an excitation point on the mastoid, and not on the jawbone [47]–[49]. Measurements were conducted the same way as described above, also with the BC devices to check the effects of airborne sound transmission. Please note that this is only part of the operation (often regarded as malfunction rather than a feature), as the main transmission should be via the skull.

Another technical parameter that can easily be determined by utilizing the setup and the equipment described above, is damping in dB over the entire frequency range in case of any headphone. This can simply be done by measuring the transmission without the headphone, followed by the same procedure with the headphones on. The reference signal in this case was the transfer characteristics of the dummy-head from the free-field to the eardrum, also called a Head-Related Transfer Function (HRTF) [50], [51]. For the measurement, a loudspeaker was set up in front of the dummy-head (frontal direction). Although HRTFs vary with direction, we only used the frontal HRTFs for damping measurements. The quotient of HRTFs with and without headphones resulted in the

Evaluation of Bone Conduction and Active-Noise-Cancellation
Headsets Based on Listening Tests in a Virtual Environment



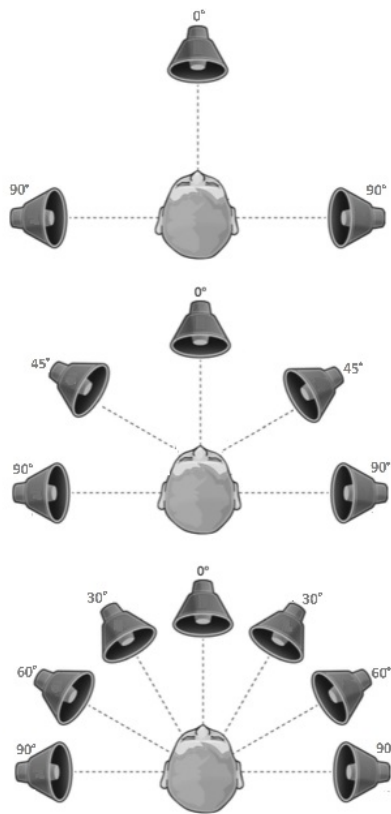Fig. 5. Damping of the ANC devices: Bose (top) and Sennheiser (bottom).

damping characteristics of the headphone (eliminating other
environmental influences, i.e., possible reflections, transfer
characteristics of the loudspeaker, etc.). In the case of BC
devices, damping is not measured as they do not cover the
outer ears. Figure 5 shows the results for the ANC types.
The blue lines are the free-field HRTFs from the frontal
direction, the black and red lines represent the transmission
if the headphones are on the ears with and without ANC. The
difference in dB gives the damping in frequency. The next
sections discuss the results of the listening test.

TABLE I
MEAN SCORES OUT OF 15 (5 DIRECTIONS)

|  | AVG | ANOVA |
|---|---|---|
| IMPULSE | 10.24 | F=1.87; Fcrit=2.70; p=0.14 |
| SPEECH | 11.92 | |
| 1KHZ SIN | 11.36 | |
| WNOISE | 11.48 | |

## B. BC

25 participants (mean age 38.2±15.1 years, 12 males and
13 females) took part in the tests. The statistical evaluation
of the comparison of the two devices showed no significant
difference between the means. Thus, the results were combined
for the purpose of simplification.

Identification of three virtual directions showed 100% ac-
curacy (9/9 points), only one subject clicked false once ac-
cidentally. This result for three directions is valid for both
BC headsets. Further evaluation was carried out by applying
five and seven directions based on ANOVA. In the case of
five directions, there was no significant difference between the
means of the four signals (F=1.87; Fcrit=2.70; p=0.14). The

TABLE II
MEAN SCORES OUT OF 21 (7 DIRECTIONS)

|  | AVG | ANOVA |
|---|---|---|
| IMPULSE | 11.12 | F=0.87; Fcrit=2.70; p=0.46 |
| SPEECH | 12.64 | |
| 1KHZ SIN | 12.20 | |
| WNOISE | 12.12 | |

same is true for 7 directions (F=0.87; Fcrit=2.70; p=0.46).
Tables I and II show the mean scores for all signal types.

## C. ANC

There were 25 participants (mean age 42.1±18.2 years,
13 males and 12 females) in the experiment. The statistical
evaluation of the comparison of the two devices showed no
significant difference between the means. Therefore, the results
were combined for the purpose of simplification.

The evaluation of three virtual directions was completely
identical to the results obtained with the use of BC. All
subjects could detect all three directions with 100% accuracy
(9/9 points), only one unintentional false click was recorded.
This result for three directions is valid for both headsets with
and without ANC activation and distractor noise.

Further evaluation was made using five and seven directions
based on ANOVA. In the case of five and seven directions,
headsets were used in three different scenarios such as:

- ANC ON - NOISE ON;
- ANC OFF - NOISE OFF;
- ANC OFF - NOISE ON.

The NOISE ON case refers to an active distractor sound.

TABLE III
MEAN SCORES OUT OF 15 (5 DIRECTIONS)

| IMPULSE | ANC ON | ANC OFF | p-value |
|---|---|---|---|
| NOISE ON | 10.6 | 10.4 | 0.864 |
| NOISE OFF | | 10.5 | |
| SPEECH | ANC ON | ANC OFF | p-value |
| NOISE ON | 12.7 | 11.6 | 0.913 |
| NOISE OFF | | 11.5 | |
| 1KHZ SIN | ANC ON | ANC OFF | p-value |
| NOISE ON | 12.1 | 12.1 | 0.769 |
| NOISE OFF | | 10.8 | |
| WNOISE | ANC ON | ANC OFF | p-value |
| NOISE ON | 13.1 | 12.7 | 0.448 |
| NOISE OFF | | 11.1 | |

TABLE IV
MEAN SCORES OUT OF 21 (7 DIRECTIONS)

| IMPULSE | ANC ON | ANC OFF | p-value |
|---|---|---|---|
| NOISE ON | 9.7 | 9.1 | 0.740 |
| NOISE OFF | | 9.4 | |
| SPEECH | ANC ON | ANC OFF | p-value |
| NOISE ON | 12.4 | 12.4 | 0.998 |
| NOISE OFF | | 11.8 | |
| 1KHZ SIN | ANC ON | ANC OFF | p-value |
| NOISE ON | 13.7 | 12.3 | 0.211 |
| NOISE OFF | | 10.5 | |
| WNOISE | ANC ON | ANC OFF | p-value |
| NOISE ON | 14.0 | 12.3 | 0.158 |
| NOISE OFF | | 12.5 | |

Tables III and IV present the results for five and seven
directions, respectively. The mean scores are shown for the

three scenarios and for all excitation signals. The p-values indicate whether there is a significant difference between ANC ON - NOISE ON and ANC OFF - NOISE ON situation.

## IV. DISCUSSION

### A. Objective evaluation

Transfer functions of the ANC devices showed typical frequency plots of high-quality headphones with the HRTF attenuation around 3 and 7 kHz, as expected. Furthermore, the effect of the ANC circuit is clearly visible causing an almost frequency independent wide band attenuation of about 10-12 dB in the case of Bose, and a negligible attenuation in the case of Sennheiser. The technical parameters of the tested devices are very similar; selection for an application can be based on subjective preference, pricing, or ergonomic considerations.

Former analyses of the transmission characteristics showed uneven frequency response of BC headphones, compared to conventional headphones or speakers [52]. Placement of the transducer on the skull is critical, especially, if a real human head is used that moves during operation. The variability in the measurements supports the subjective findings about the sound quality to be influenced by the placement. There is a significant airborne sound transmission from about 200-400 Hz up to 10 kHz in case of an open ear canal.

Increased damping can be observed between 100 Hz and 1000 Hz, using the ANC mode by both ANC devices, up to 20 dB. Above 1 kHz, there is no significant difference with or without ANC; however, the headphone damping becomes significant. We can conclude that the missing natural damping of the headphone below 1 kHz can be extended to about 100 Hz by activating ANC.

For BC devices, damping of airborne sound could be assured if the ear canal entrances were plugged. However, the main goal of this equipment is to leave the ear canal open as the primary operation. In a former measurement, the lowest threshold (i.e., maximum sensitivity) in both open and plugged conditions occurred around 1170 Hz-1370 Hz and within the frequency range of speech. At frequencies below or above these frequency values, greater intensity was required to detect the sound output by the bonephones [23]. Plugging the ears would lower the threshold up to about 2000 Hz.

### B. Subjective evaluation

The main goals of the experiment were to compare four individual signal types in three different virtual speaker setups, testing the effect of ANC if a distractor sound is present, and to decide whether ANC and BC headsets have benefits in contrast to regular headphones.

The AudioBone headset was also used in a listening test where standard speaker setups were contrasted with BC audio. A speaker constellation of 90°, 45° and 0° was installed, and stationary and moving sources were emulated. Results showed the best accuracy with a physical speaker array and stationary sound, but there was no difference between the speaker array and the bone conduction device for sounds that

were moving [31]. Another study used the Aftershokz Sportz3 [34]. Here the goal was to introduce vertical localization cues to BC playback. It was reported that this could result in decreased localization accuracy in the horizontal plane. Furthermore, there was also a significant "compression" in the area directly in front of the observer in elevation. Our measurements indicate a similar effect horizontally in the case of seven directions, as subjects reported 30°-60° separation to be more difficult than separating 60° and 90°.

Employing the BC devices, participants reported upon their subjective feelings related to the signals. Impulse and noise were the most annoying and the most difficult ones to localize, and speech was the best/easiest. Indeed, results with impulse were behind the results of the others. On the other hand, the subjective feelings on white noise were not supported by the results. All but one participant mentioned this signal to be the most difficult one to localize; nevertheless, their performance was better than with other signals. As expected, it was a common observation that having seven instead of five directions increases difficulty. Seven directions were reported to be inappropriate and "too many", while three directions were "too easy". We can hypothesize that different results would be measured if the three directions were not 90° apart, but closer to one another. Subjects seemed to be connecting the unpleasant sounds to the difficulty of localization (i.e., noise and impulse are disturbing, therefore, more challenging to localize in contrast to speech). They were, however, not related. It was also often mentioned that it is harder to distinguish between 30° and 60° than between 60° and 90° in the 7-direction scenario. Furthermore, results may be biased by the actual order of the randomized signal presentation; directions can be compared more reliable if they are neighbors. For example, directions 30° and 60° from the same side can be easier separated if they are presented after one another.

In case of BC, the placement and adjustment of the transducer on the jawbone affect the observed sound quality. Displacement during the listening test due to small head movements, swallowing, moving the jaws can result in a drop of sound quality, and losing the symmetry between the left- and the right-hand sides. It is common that subjects frequently replace and adjust the position of the transducers. Furthermore, the AfterShokz has a belt that cannot be adjusted to the head size and the fixed transducer position is not optimal for all subjects.

The effects of learning and experience were not measured directly. It was, however, self-reported by the subjects that they had gotten better at the task after completing several sessions. This also had a motivational effect; some users tried to become better and "break the record". We could find evidence from the literature about the effect of training in listening tests. Although not tested directly, we can support the findings that the localization performance may be increased by repeated sessions [53]–[56].

Similarly, the experiment was not designed to test the dependence between the results and age. Elderly subjects (45+) were underrepresented in the sample. Although a statistical

analysis showed no difference between the groups above and below 45, a correct sample of subjects with an equal number of participants in age groups may result in a different outcome.

There is no need to compare anything in the case of three directions, as subjects delivered a perfect accuracy in all cases. The mean scores and the corresponding ANOVA results for five and seven directions can be seen in Tables I and II. There is no significant difference among the four excitation signals.

Based on Table V, the average scores vary from 68.3% to 79.5% for five directions and from 52.9% to 60.2% for seven directions, respectively. As expected, there is a significant difference (F=39.35; Fcrit=5.98; p=0.0007).

TABLE V
COMPARISON OF AVERAGE SCORES FOR EACH SIGNAL (BC TYPES ONLY
FOR 5 AND 7 DIRECTIONS)

| | 5 directions AVG SCORE OUT OF 15 | % | 7 directions AVG SCORE OUT Of 21 | % |
|---|---|---|---|---|
| IMPULSE | 10.24 | 68.3 | 11.12 | 52.9 |
| SPEECH | 11.92 | 79.5 | 12.64 | 60.2 |
| 1 KHZ SIN | 11.36 | 75.7 | 12.20 | 58.1 |
| WNOISE | 11.48 | 76.5 | 12.12 | 57.7 |

For the ANC headsets, the mean scores and the corresponding ANOVA results for five and seven directions are depicted in Tables III and IV. There is no significant difference among the four excitation signals with and without ANC and external noise.

Based on Table VI the average scores vary from 65.6% to 76% for five directions and from 45.6% to 59.2% for seven directions, respectively. As expected, there is a significant difference (F=24.42; Fcrit=5.98; p=0.0026).

TABLE VI
COMPARISON OF AVERAGE SCORES FOR EACH SIGNAL (ANC TYPES ONLY
FOR 5 AND 7 DIRECTIONS)

| | 5 directions AVG SCORE OUT OF 15 | % | 7 directions AVG SCORE OUT Of 21 | % |
|---|---|---|---|---|
| IMPULSE | 9.84 | 65.6 | 9.58 | 45.6 |
| SPEECH | 10.92 | 72.8 | 11.87 | 56.5 |
| 1 KHZ SIN | 11.28 | 75.2 | 11.28 | 53.7 |
| WNOISE | 11.40 | 76.0 | 12.43 | 59.2 |

One may conclude that the impulse signal is the most difficult one to use, even though the difference is statistically not significant at the 0.05 level. Indeed, it may be related to the fact that the first signal was always the impulse; therefore, the subjects' performance improved based on user experience.

As expected, increasing the number of source locations from five to seven drops drastically the number of correct scores. An exhaustive paired statistical analysis across all measurements for all four signals between five and seven directions in the ANC ON-NOISE OFF (p=0.009, p=0.033, p=0.01, p=0.035); ANC ON-NOISE ON (p=0.035, p=0.003, p=0.045, p=0.011); and ANC OFF-NOISE ON (p=0.002, p=0.004, p=0.003, p=0.001) scenarios revealed statistically significant differences in each case. For the drastic drop of correct judgements (from 14 to 26%) is the introduction of

two more directions responsible only, while the type of the excitation signal, ANC activation and the presence or absence of the external distractor sound are insignificant.

Comparing BC and ANC headsets is an interesting question. The combined results of the two BC and ANC devices respectively, showed no significant difference neither in case of five directions (F=1.44; Fcrit=3.49; p=0.28) nor in case of seven directions (F=0.67; Fcrit=3.49; p=0.59). Looking at the means over all signal types, this result could be expected. Average scores vary from 10.4 to 13.1 out of 15; and from 9.1 to 14 out of 21. We concluded that in this experiment both BC devices and both ANC headsets perform equally. Furthermore, if we consider ANC OFF-NOISE OFF situation as a reference traditional headphone case (all ANC devices act as a regular headphone if the battery is low or ANC is off), paired t-tests support that standard headphone modes do not outperform BC devices and will not be outperformed by an activated ANC mode.

In this experiment, no external distractor sounds were used in the case of BC headsets. The ANC effect was tested with simulated traffic noise. The additional damping of the ANC circuit did not affect the localization performance at all. We may expect different outcomes in case of higher SPL of the distractor sound. The applied 80 dB SPL was reduced by the damping effect to about 20 dB; thus, masking effect of the noise was not detected. The masking effect of the distractor sounds is a widely investigated area [57]–[61]. In a reverse condition where the distractor sound came from the BC set, listeners had to localize real sound sources (loudspeakers) located around them, with and without distractor sounds. Participants had greater localization error in the distractor-present conditions, especially in case of multiple distractors, regardless of whether or not participants ignored distractors [62].

Although not included in the final experiment, other models were also used in preliminary tests during the design-phase of the experiment. The Sony WH-1000XM3 model is a supraaural ANC headset, the Bose QC20 is an in-ear phone version of an ANC type device, and we also used the wireless AfterShokz BC model. Each device seemed to be applicable, they differed generally in comfort, size, weight and ergonomics.

Note that the task in this experiment was to detect and identify a limited number of possible sound sources. This is different than a so-called absolute localization task where the possible sound direction is "continuous" in space. It is easier to select a sound source from three or five discrete positions than pointing to one perceived location. If having a limited, discrete number of potential source locations, subjects often do not localize at all. Nevertheless, they make their judgements based on the available information, which - in this case - was the visual representations of virtual loudspeakers. Subjects were not asked about externalization directly, but they did not mention any in-the-head localization problems during their selections. Even if they had it, they were not aware of it and not influenced by this phenomenon, neither with ANC nor with BC devices. Applying complex binaural rendering

with HRTFs, individual parameter settings, equalization of headphones, etc., is not needed for simple applications. Using binaural spatialization in the horizontal plane with a BC headset, the minimum discernable angular difference between two successively presented sound sources was found to be around 10° and above [34]. However, accuracy was poorer than with headphone based reproduction. Localization errors between 30°–35° were measured in front, back, and sides. HRTFs were applied also in an experiment targeting the usability of BC headsets for VR applications. Azimuth accuracy was about the same as it was for a regular headphone, while for the same result in elevation accuracy, high frequency components were needed [63]. Both experiments applied binaural spatialization to determine the absolute localization accuracy.

## V. CONCLUSIONS

A virtual listening test was installed in a VR scenario using three, five and seven source directions and four different excitation signals. 25 untrained subjects evaluated two bone conduction and another 25 participants evaluated two ANC-equipped headphones. The results showed that there is no statistically significant difference between the test signals. Furthermore, applying an 80 dB external distractor sound did not influence localization judgements with or without active noise cancellation. Applying stereo panning, three directions (F/L/R) can be used any time without errors. On the other hand, there is a significant difference between the mean scores of the 5-direction and the 7-direction scenarios in all cases. Five directions (in 45-degree spacing in the frontal hemisphere) can be used with 65-75% accuracy and it is recommended for real applications where errors are not critical. Based on the measurement results and the participants' subjective evaluations, seven directions are too demanding and inaccurate for applications. All ANC and BC devices performed equally in accuracy in a given task.

Future directions for research include recruiting visually impaired users, using different external or internal distractor sounds for BC and ANC, testing looped playback and a closed ear canal. Furthermore, rearranged virtual speaker setups (3-5 directions), testing in real life scenarios, and modification of the virtual environment to an explorable dynamic 3D scenario could provide additional information about the usability of the devices.

## REFERENCES

[1] J. C. Middlebrooks and D. M. Green, "Sound localization by human listeners," *Annu. Rev. Psychol*, vol. 42, pp. 135–59, 1991. DOI: 10.1146/annurev.ps.42.020191.001031

[2] M. Vorländer, *Auralization*. Heidelberg: Springer, 2020.

[3] W. M. Hartmann, *Localization and Lateralization of Sound*. Cham: Springer International Publishing, 2021, pp. 9–45.

[4] G. Plenge, "On the differences between localization and lateralization," *The Journal of the Acoustical Society of America*, vol. 56, no. 3, pp. 944–951, 1974. DOI: 10.1121/1.1903353

[5] B. N. Walker and J. Lindsay, "Using virtual environments to prototype auditory navigation displays," *Ass. Technology*, vol. 17, no. 1, pp. 72–81, 2005. DOI: 10.1080/10400435.2005.10132097

[6] B. Kapralos, M. Jenkin, and E. Milios, "Virtual audio systems," *Presence: Teleoperators and Virtual Environments*, vol. 17, no. 6, pp. 527–549, 2008. DOI: 10.1162/pres.17.6.527

[7] T. Kuppanda, N. Degara, D. Worrall, B. Thoshkahna, and M. Müller, "Virtual reality platform for sonification evaluation," in *Proceedings of ICAD'15*, Graz, Austria. GeorgiaTech, 2015, pp. 117–124.

[8] G. C. Stecker, "Using virtual reality to assess auditory performance," *The Hearing Journal*, vol. 72, no. 6, p. 20, 2019. DOI: 10.1097/01.hj.0000558464.75151.52

[9] V. Hohmann, R. Paluch, M. Krueger, M. Meis, and G. Grimm, "The virtual reality lab: realization and application of virtual sound environments," *Ear and Hearing*, vol. 41, no. Suppl 1, p. 31S, 2020. DOI: 10.1097/AUD.0000000000000945

[10] J. Tonndorf, "Bone conduction." Berlin: Springer, 1976, pp. 37–84.

[11] S. Stenfelt, "Acoustic and physiologic aspects of bone conduction hearing," *Implantable bone conduction hearing aids*, vol. 71, pp. 10–21, 2011. DOI: 10.1159/000323574

[12] C. N. Hansen, *Understanding active noise cancellation*. Oxford: CRC Press, 1999.

[13] C.-Y. Chang and S.-T. Li, "Active noise control in headsets by using a low-cost microcontroller," *IEEE Transactions on industrial electronics*, vol. 58, no. 5, pp. 1936–1942, 2010. DOI: 10.1109/TIE.2010.2058071

[14] S. Spagnol, G. Wersényi, M. Bujacz, O. Bălan, M. Herrera Martínez, A. Moldoveanu, and R. Unnthorsson, "Current use and future perspectives of spatial audio technologies in electronic travel aids," *Wireless Communications and Mobile Computing*, vol. 2018, 2018. DOI: 10.1155/2018/3918284

[15] J. C. Lock, I. D. Gilchrist, G. Cielniak, and N. Bellotto, "Bone-conduction audio interface to guide people with visual impairments," in *International Conference on Smart City and Informatization*. Springer, 2019., pp. 542–553, DOI: 10.1007/978-981-15-1301-5-43.

[16] J. A. MacDonald, P. P. Henry, and T. R. Letowski, "Spatial audio through a bone conduction interface: Audición espacial a través de una interfase de conducción ósea," *International journal of audiology*, vol. 45, no. 10, pp. 595–599, 2006. DOI: 10.1080/14992020600876519

[17] B. F. Katz and L. Picinali, "Spatial audio applied to research with the blind," *Advances in sound localization*, pp. 225–250, 2011. DOI: 10.5772/15206

[18] G. A. Studebaker, "Placement of vibrator in bone-conduction testing," *Journal of speech and hearing research*, vol. 5, no. 4, pp. 321–331, 1962. DOI: 10.1044/jshr.0504.321

[19] M. McBride, T. Letowski, and P. Tran, "Bone conduction reception: Head sensitivity mapping," *Ergonomics*, vol. 51, no. 5, pp. 702–718, 2008. DOI: 10.1080/00140130701747509

[20] R. H. Margolis, R. H. Eikelboom, C. Johnson, S. M. Ginter, D. W. Swanepoel, and B. C. Moore, "False air-bone gaps at 4 khz in listeners with normal hearing and sensorineural hearing loss," *International journal of audiology*, vol. 52, no. 8, pp. 526–532, 2013. DOI: 10.3109/14992027.2013.792437

[21] B. Håkansson, A. Tjellström, and U. Rosenhall, "Hearing thresholds with direct bone conduction versus conventional bone conduction," *Scandinavian audiology*, vol. 13, no. 1, pp. 3–13, 1984. DOI: 10.3109/01050398409076252

[22] N. A. Watson, "Limits of audition for bone conduction," *The Journal of the Acoustical Society of America*, vol. 9, no. 4, pp. 294–300, 1938. DOI: 10.1121/1.1915936

[23] B. N. Walker and R. M. Stanley, "Thresholds of audibility for bone-conduction headsets," in *Proceedings of ICAD'05, Limerick, Ireland. Georgia Institute of Technology*, 2005, pp. 218–222.

[24] S. Stenfelt, "Transcranial attenuation of bone-conducted sound when stimulation is at the mastoid and at the bone conduction hearing aid position," *Otology & neurotology*, vol. 33, no. 2, pp. 105–114, 2012. DOI: 10.1097/MAO.0b013e31823e28ab

[25] S. Stenfelt and M. Zeitooni, "Binaural hearing ability with mastoid applied bilateral bone conduction stimulation in normal hearing subjects," *The Journal of the Acoustical Society of America*, vol. 134, no. 1, pp. 481–493, 2013. DOI: 10.1121/1.4807637

[26] M. Zeitooni, E. Mäki-Torkko, and S. Stenfelt, "Binaural hearing ability with bilateral bone conduction stimulation in subjects with normal hearing: Implications for bone conduction hearing aids," *Ear and hearing*, vol. 37, no. 6, pp. 690–702, 2016. DOI: 10.1097/AUD.0000000000000336

[27] D. Schonstein, L. Ferré, and B. F. Katz, "Comparison of headphones and equalization for virtual auditory source localization," *Journal of the Acoustical Society of America*, vol. 123, no. 5, p. 3724, 2008. DOI: 10.1121/1.2935199

[28] A. F. Jahn and J. Tonndorf, "Lateralization of bone-conducted sounds," *American journal of otolaryngology*, vol. 3, no. 2, pp. 133–140, 1982. **DOI**: 10.1016/S0196-0709(82)80044-6

[29] K. Kaga, M. Setou, and M. Nakamura, "Bone-conducted sound lateralization of interaural time difference and interaural intensity difference in children and a young adult with bilateral microtia and atresia of the ears," *Acta oto-laryngologica*, vol. 121, no. 2, pp. 274–277, 2001. **DOI**: 10.1080/000164801300043820

[30] B. N. Walker, R. M. Stanley, N. Iyer, B. D. Simpson, and D. S. Brungart, "Evaluation of bone-conduction headsets for use in multitalker communication environments," in *Proceedings of the human factors and ergonomics society annual meeting*, vol. 49, no. 17. SAGE Publications Sage CA: Los Angeles, CA, 2005. pp. 1615–1619, **DOI**: 10.1177/154193120504901725.

[31] R. W. Lindeman, H. Noma, and P. G. De Barros, "Hear-through and mic-through augmented reality: Using bone conduction to display spatialized audio," in *2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality*. IEEE, 2007. pp. 173–176, **DOI**: 10.1109/ISMAR.2007.4538843.

[32] R. M. Stanley, "Measurement and validation of bone-conduction adjustment functions in virtual 3d audio displays," Ph.D. dissertation, Georgia Institute of Technology, 2009.

[33] G. Wersényi, "Localization in a head-related transfer function-based virtual audio synthesis using additional high-pass and low-pass filtering of sound sources," *Acoustical science and technology*, vol. 28, no. 4, pp. 244–250, 2007. **DOI**: 10.1250/ast.28.244

[34] A. Barde, G. Lee, W. S. Helton, and M. Billinghurst, "Binaural spatialisation over a bone conduction headset: Elevation perception," in *Proceedings of ICAD'16*, Canberra, Australia. International Community on Auditory Display, 2016, pp. 1–4.

[35] S. M. Kuo, K. Kuo, and W. S. Gan, "Active noise control: Open problems and challenges," in *The 2010 International conference on green circuits and systems*. IEEE, 2010., pp. 164–169, **DOI**: 10.1109/ICGCS.2010.5543076.

[36] M. Guldenschuh, A. Sontacchi, M. Perkmann, and M. Opitz, "Assessment of active noise cancelling headphones," in *2012 IEEE Second International Conference on Consumer Electronics-Berlin (ICCE-Berlin)*. IEEE, 2012., pp. 299–303, **DOI**: 10.1109/ICCE-Berlin.2012.6336504.

[37] S. Liebich, J. Fabry, P. Jax, and P. Vary, "Signal processing challenges for active noise cancellation headphones," in *Speech Communication; 13th ITG-Symposium*. VDE, 2018, pp. 1–5.

[38] "MaxWhere 3D Platform," https://www.maxwhere.com//, 2022.

[39] B. Berki, "Better memory performance for images in maxwhere 3d vr space than in website," in *2018 9th IEEE International Conference on Cognitive Infocommunications (CogInfoCom)*, 2018., pp. 281–284, **DOI**: 10.1109/CogInfoCom.2018.8639956.

[40] "AfterShokz," https://shokz.com///, 2022.

[41] "AudioBone," https://www.goldendance.co.jp/English/product/p-ab01.html//, 2022.

[42] "Sennheiser," https://en-au.sennheiser.com/pxc450//, 2022.

[43] "Bose QuietComfort 25," {https://www.bose.com//}, 2022.

[44] H. Møller, D. Hammershøi, C. B. Jensen, and M. F. Sørensen, "Transfer characteristics of headphones measured on human ears," *Journal of the Audio Engineering Society*, vol. 43, no. 4, pp. 203–217, 1995.

[45] A. Kulkarni and H. S. Colburn, "Variability in the characterization of the headphone transfer-function," *The Journal of the Acoustical Society of America*, vol. 107, no. 2, pp. 1071–1074, 2000. **DOI**: 10.1121/1.428571

[46] G. Wersényi, "Comparison of Transfer Functions of Open Ear Canal Headsets Measured on a Dummy-Head and a Human Head," in *Forum Acusticum 23*. EAA, 2023, p. 5.

[47] L. A. Wilber and V. Goodhill, "Real ear versus artificial mastoid methods of calibration of bone-conduction vibrators," *Journal of Speech and Hearing Research*, vol. 10, no. 3, pp. 405–416, 1967. **DOI**: 10.1044/jshr.1003.405

[48] R. H. Margolis and S. M. Stiepan, "Acoustic method for calibration of audiometric bone vibrators," *The Journal of the Acoustical Society of America*, vol. 131, no. 2, pp. 1221–1225, 2012. **DOI**: 10.1121/1.3675007

[49] B. Karaböce, "Force sensitivity determination of artificial mastoid used for audiometric bone-conduction measurements," in *2020 IEEE International Symposium on Medical Measurements and Applications (MeMeA)*. IEEE, 2020., pp. 1–5. **DOI**: 10.1109/MeMeA49120.2020.9137333.

[50] S. Li and J. Peissig, "Measurement of head-related transfer functions: A review," *Applied Sciences*, vol. 10, no. 14, p. 5014, 2020. **DOI**: 10.3390/app10145014

[51] C. I. Cheng and G. H. Wakefield, "Introduction to head-related transfer functions (hrtfs): Representations of hrtfs in time, frequency, and space," in *AEA Convention 107*. Audio Engineering Society, 1999.

[52] T. M. Voong and M. Oehler, "Auditory spatial perception using bone conduction headphones along with fitted head related transfer functions," in *2019 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. IEEE, 2019., pp. 1211–1212, **DOI**: 10.1109/VR.2019.8798218.

[53] M. Ohuchi, Y. Iwaya, Y. Suzuki, and T. Munekata, "Training effect of a virtual auditory game on sound localization ability of the visually impaired," in *Proc. of ICAD'05*, Limerick, Ireland, 2005, pp. 283–286.

[54] P. Majdak, M. J. Goupell, and B. Laback, "3-d localization of virtual sound sources: Effects of visual environment, pointing method, and training," *Attention, perception, & psychophysics*, vol. 72, no. 2, pp. 454–469, 2010. **DOI**: 10.3758/APP.72.2.454

[55] M. A. Steadman, C. Kim, J.-H. Lestang, D. F. Goodman, and L. Picinali, "Short-term effects of sound localization training in virtual reality," *Scientific Reports*, vol. 9, no. 1, pp. 1–17, 2019.

[56] S. Bech, "Selection and training of subjects for listening tests on sound-reproducing equipment," *Journal of the Audio Engineering Society*, vol. 40, no. 7/8, pp. 590–610, 1992.

[57] M. D. Good, R. H. Gilkey, and J. M. Ball, "The relation between detection in noise and localization in noise in the free field," *Binaural and spatial hearing in real and virtual environments*, pp. 349–376, 1997.

[58] E. A. Macpherson and J. C. Middlebrooks, "Localization of brief sounds: effects of level and background noise," *The Journal of the Acoustical Society of America*, vol. 108, no. 4, pp. 1834–1849, 2000. **DOI**: 10.1121/1.1310196

[59] D. S. Brungart, B. D. Simpson, and A. J. Kordik, "Localization in the presence of multiple simultaneous sounds," *Acta Acustica united with Acustica*, vol. 91, no. 3, pp. 471–479, 2005.

[60] M. F. Mueller, A. Kegel, S. M. Schimmel, N. Dillier, and M. Hofbauer, "Localization of virtual sound sources with bilateral hearing aids in real- istic acoustical scenes," *The Journal of the Acoustical Society of America*, vol. 131, no. 6, pp. 4732–4742, 2012. **DOI**: 10.1121/1.4705292

[61] K. Lopatka, J. Kotus, and A. Czyzewski, "Detection, classification and localization of acoustic events in the presence of background noise for acoustic surveillance of hazardous situations," *Multimedia Tools and Applications*, vol. 75, no. 17, pp. 10 407–10 439, 2016. **DOI**: 10.1007/s11042-015-3105-4

[62] K. R. May and B. N. Walker, "The effects of distractor sounds presented through bone conduction headphones on the localization of critical environmental sounds," *Applied ergonomics*, vol. 61, pp. 144–158, 2017. **DOI**: 10.1016/j.apergo.2017.01.009

[63] M. V. Tray, M. do VM da Costa, M. Regener, C. Reuter, and M. Oehler, "Suitability of bone conduction headphones for virtual acoustic environments," in *AES 2022 Int. Audio for Virtual and Augmented Reality Conference*. Audio Engineering Society, 2022.

**György Wersényi** was born in 1975 in Győr, Hungary. He received his MSc degree in electrical engineering from the Technical University of Budapest in 1998 and PhD degree from the Brandenburg Technical University in Cottbus, Germany. Since 2002 he has been member of the Department of Telecommunications at the Széchenyi István University in Győr. From 2020 to 2022 he was the dean of Faculty of Mechanical Engineering, Informatics and Electrical Engineering, as well as the scientific president of the Digital Development Center at the university. Currently, he is a full professor, member of the European Acoustics Association (EAA) and the Audio Engineering Society (AES). His research focus is on acoustic measurements, virtual and augmented reality solutions, sonification, cognitive infocommunications, and assistive technologies.

# Quantum Genetic Algorithm for Highly Constrained Optimization Problems

Abdulbasit M. A. Sabaawi[1,2], Mohammed R. Almasaoodi[1,3], Sara El Gaily[1], and Sándor Imre[1]

*Abstract*—**Quantum computing appears as an alternative solution for solving computationally intractable problems. This paper presents a new constrained quantum genetic algorithm designed specifically for identifying the extreme value of a highly constrained optimization problem, where the search space size _database is massive and unsorted_ cannot be handled by the currently available classical or quantum processor, called the highly constrained quantum genetic algorithm (HCQGA). To validate the efficiency of the suggested quantum method, maximizing the energy efficiency with respect to the target user bit rate of an uplink multi-cell massive multiple-input and multiple- output (MIMO) system is considered as an application. Simulation results demonstrate that the proposed HCQGA converges rapidly to the optimum solution compared with its classical benchmark.**

*Index Terms*—**genetic algorithm, quantum computing, quantum extreme value searching algorithm, blind quantum computing**

## I. INTRODUCTION

Quantum computing technology provides efficient solutions to handle intractable computational problems that classical computers are unable to handle. This is achieved by exploiting the fundamentals of quantum nature, such as quantum superposition and quantum entanglement [1], [2]. Different communities and research organizations have used quantum computing with the promise of solving various classes of computational problems in many areas, including security and cryptography [3]–[5], networks [6]–[9], space communication [10], [11], and many other optimization fields. Quantum computing guarantees exponential speed, short running time, and high accuracy [12].

Optimal decisions with respect to certain constraints are crucial in various applications in a variety of areas, including physics, communications, and computer science. For instance, in computer networking, routing algorithms must factor in network congestion, latency, and reliability to find the best data path between devices. Similarly, in wireless communications, mobile providers must efficiently allocate radio spectrum to ensure fast and reliable connections without interference or congestion [13].

There are two forms of constrained optimization problems: the first type refers to non-highly constrained; it means that the constraints are relatively lenient, resulting in a larger feasible set of solutions. This can make satisfying the constraints simpler, and it may also increase the number of potential solutions. The second type consists of highly constrained problems, in which constraints are stringent and limit the feasible set of solutions, which can lead finding a solution more difficult but can also lead to the optimal and desired solution. On the other hand, in real-world optimization problems, the search space, or size of the candidate solutions for a given optimization problem, is exceedingly vast, rendering it beyond the computational capacity of both classical and quantum computers. Additionally, the search space remains unsorted, further complicating the task at hand.

The challenges presented by the highly constrained optimization problems, as well as the search space size that surpasses the computational capacities of existing classical and quantum machines, have motivated us to develop a novel quantum genetic algorithm called the highly constrained quantum genetic algorithm (HCQGA). An uplink multi-cell massive multiple-input and multiple-output (MIMO) application is utilized as a representative example for testing and validating the efficiency of the HCQGA. Note that this paper is an extension of the preceding research publications [14]–[19].

Recently, researchers have been significantly more interested in establishing various quantum genetic algorithms to solve optimization problems. The authors in [20] proposed a QGA based scheduling algorithm for heterogeneous platforms, the algorithm improves task scheduling efficiency and reduces the computational cost, but there are also communication-intensive tasks in signal processing. they reduced it by task clustering or task duplication. In [21], the authors presented a QGA with simultaneous quantum crossover on all chromosomes. They utilized two identical copies of a superposition for qubit relabeling, a complexity analysis implies that a quadratic speedup over its classical counterpart is attained in each generation's dominant factor of the run time. In [22], the authors improved the initial population stage of the genetic algorithm by employing the quantum counting algorithm to detect the number of unsuitable chromosomes in the population, the goal of this paper is to exploit a quantum algorithm faster than the classical counterparts on the optimizing performance. In [23] introduced an immigration technique that enhances the QGA by considering the most optimal qubit string in quantum chromosomes. Randomly

[1] Department of Networked Systems and Services, Faculty of Electrical Engineering and Informatics, Budapest University of Technology and Economics, Budapest, Hungary.
[2] College of Electronics Engineering, Ninevah University, Mosul, Iraq.
[3] Kerbala University, Kerbala, Iraq.
E-mail: {abdulbasit, almasaoodi, elgaily, imre}@hit.bme.hu

transferring this highly fit qubit string to the next generation's chromosomes leads to improved mixing and overall algorithm performance.

To the best of our knowledge, there is currently no QGA available that addresses the aforementioned problem statement. This is primarily due to the limitations of most QGAs, which can only be executed on a universal quantum computer with a constrained qubit size. Consequently, in this research, we did not explore a comparison between the existing QGAs and our proposed HCQGA.

Massive Multiple-Input Multiple-Output (MIMO) is an advanced form of MIMO technology, featuring a large number of antennas at base stations to concurrently serve multiple users. By leveraging spatial diversity and array gain, it enhances system capacity, spectral efficiency, and energy efficiency. Massive MIMO outperforms in areas with significant population density due to its superior interference management, signal quality, and coverage extension. This technology is an integral component of 5G and 6G networks, delivering higher data rates and enhanced wireless communication performance to satisfy the expanding demands of modern connectivity [24]–[26]. Maximizing energy efficiency (EE) holds significant importance due to its potential for mitigating environmental impact, conserving finite resources, and yielding substantial cost savings [27], [28].

The EE performance of massive MIMO has been explored in the following studies: In [29], the authors proposed an energy-efficient low-complexity algorithm, using Newton's methods and Lagrange's decomposition for optimal power allocation and user association in massive MIMO. In [30], an adaptive power allocation algorithm based on particle swarm optimization for enhancing energy efficiency in uplink multi-cell, multi-user massive MIMO systems was proposed. However, it does not consider the impact of Channel State Information (CSI) on this optimization in uplink multi-cell massive MIMO systems. The authors in [31] suggested a sub-optimal Dinkelbach-like algorithm to enhance energy efficiency in the downlink of a multi-cell multi-user massive MIMO system, under interference limitations and maximum transmit power constraints.

The subsequent sections of this paper are organized as follows: Section II provides a comprehensive overview of the fundamental background required for the development of the HCQGA, such as the quantum extreme value searching algorithm (QEVSA), blind quantum computing (BQC), the unconstrained classical genetic algorithm (UCGA), and the unconstrained quantum genetic algorithm (UQGA). Section III introduces and extends the UQGA version into HCQGA. In Section IV, the HCQGA is applied in an uplink multi-cell massive MIMO system in order to maximize energy efficiency with respect to the desired bit rate of users. To validate the efficiency of the HCQGA, simulation experiments have been conducted in Section V. Finally, Section VI concludes the manuscript.

## II. METHODS

To better comprehend the novelty of the developed HCQGA, it is essential to introduce the preliminary strategies that form its foundation. These strategies include the QEVSA, BQC, UCGA, and UQGA. By merging these strategies, the HCQGA presents a unique and innovative approach that addresses the highly constrained optimization problems, which entail search space sizes surpassing the computational capacities of existing classical and quantum machines.

### A. Quantum Extreme Value Searching Algorithm

The QEVSA is designed to search for the extreme value of an unconstrained cost function or an unsorted database. It combines two distinct concepts: the classical binary searching algorithm (BSA), which utilized for identifying a specific element/item within a sorted database [32], [33], and quantum existence testing (QET) [34], which is a special case of the quantum counting algorithm. When the QET detects the presence of the searched target in the unsorted database, it returns the answer "YES" otherwise, "NO" [35]. The QEVSA is given in detail as follows:

| | QEVSA |
|---|---|
| 1 | We begin with step $h = 0$ : $F_{min\,1} = F_{min\,0}$, $F_{max\,1} = F_{max\,0}$, and $\Delta F = F_{max\,0} - F_{min\,0}$ |
| 2 | $h = h + 1$ |
| 3 | $F_{med\,h} = F_{min\,h} + \left[\dfrac{F_{max\,h} - F_{min\,h}}{2}\right]$ |
| 4 | $flag = QET\,(F_{med\,h})$: <br><br> • if $flag = Yes$, then $F_{max\,h+1} = F_{med\,h}$   $F_{min\,h+1} = F_{min\,h}$ <br><br> • Else $F_{max\,h+1} = F_{max\,h}$ , $F_{min\,h+1} = F_{med\,h}$ |
| 5 | If $h < log_2\,(T)$, then go to 2, else stop and $y_{opt} = F_{med\,h}$ |

The parameter $T$ denotes the overall number of steps required to execute the BSA embedded in the QEVSA. While the parameter $F$ denoted the goal function, such that $F_{med\,h}$ denotes the cost function. The computational complexity of the QEVSA is $O\left(log_2(T)\,log_2{}^3\left(\sqrt{N}\right)\right)$ steps, where $N$ denotes the overall search space.

### B. Blind Quantum Computing

Blind quantum computing aims to delegate computations to untrusted quantum computes/quantum slave nodes. The BQC structure consists of one quantum server connected to several quantum computing nodes via the internet, i.e., the quantum computing nodes are accessible remotely thanks to the development of quantum communication over optical fiber.

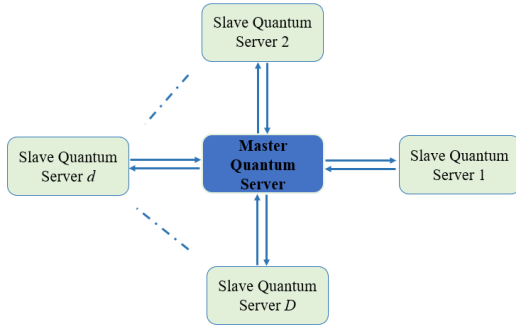Figure 1 illustrates the blind quantum computing architecture.



Fig. 1. Blind quantum computing architecture.

The BQC also ensures the privacy and security of data while enabling computations to be performed directly on encrypted information. The working mechanism of the QBC involves several steps, including encryption, computation on encrypted data, and decryption, to enable the processing of data while maintaining its privacy and confidentiality. Moreover, it offers significant benefits due to processing sensitive data without the need for decoding, thereby reducing the risks associated with exposure to potential threats [36].

*C. Unconstrained Classical Genetic Algorithm*

The UQGA is a metaheuristic strategy. It is commonly applied to solving optimization problems. There are several versions of genetic algorithms that commonly share similar working steps. The UQGA starts by randomly generating an initial population (subset from the overall set of candidate solutions). A possible candidate solution is called the chromosome. The quality of the chromosomes in a population is evaluated using a fitness function, i.e., the chromosome is assigned a fitness function value that indicates whether it approaches the optimal result or not. The present chromosomes of the population are subsequently evolved through repeated iterations, called generations (in each generation, the chromosomes undergo an evaluation process). Next, half of the chromosomes of the population with the highest fitness value are chosen for reproducing a new population (this set is named the parent set). This phase is called selection stage. Then, the selected half of the chromosomes undergo the process of crossover and mutation. The resulting new set (the second produced half of chromosomes in the population) is referred to as the offspring set. The next population is generated by merging both the parent and offspring set. The repetition of these steps for every new population aims to attain the best possible chromosome. The UCGA presents in detail as follows:

| UCGA | |
|---|---|
| 1 | Start with $step = 0$ and initialize the population size $b$. |
| 2 | Generate random population $P_s$. |
| 3 | Compute the fitness of $P_s$ and execute a sorting algorithm to extract the parent set. |
| | Crossover and mutation operations. |
| | Produce $P_{s+1}$. |
| 4 | If the global solution $J_{opt}$ is found, then stop, else go to 3. |

*D. Unconstrained Quantum Genetic Algorithm*

The UQGA combines the power of the QEVSA, QBC, and UCGA to find the extreme value in a large and unsorted database, in which no classical or quantum computer can handle the search [16].

The UQGA replaces the random initialization stage and the classical selection stage with their quantum counterparts versions:

- The quantum initialization stage: This stage increases the quality of chromosomes and has a significant impact on the convergence speed of the UQGA to the optimal result to produce high quality chromosomes in a population. The quantum server selects a set of chromosomes (regions) randomly with the same sizes denoted by $R$. Next, the quantum server assigns every region to a quantum slave node. Then, the QEVSA is applied in every region in order to extract the initial population. The computational complexity of this stage is $O\left(log_2(T)log_2{}^3\left(\sqrt{R}\right)\right)$ steps.

- The quantum selection stage: The procedure consists of applying $\frac{U}{2}$ times the QEVSA instead of applying the classical sorting algorithm. The computational complexity of this stage is $O\left(\frac{U}{2}.log_2(T)log_2{}^3\left(\sqrt{U}\right)\right)$ steps.

| UQGA | |
|---|---|
| 1 | Start with $step = 0$. Set the population and sub-database sizes, respectively $U$ and $G$. |
| 2 | Generate the $U$ regions. |
| 3 | Run the QEVSA for every sub-database to create population $P_{st}$. |
| 4 | Select the first half of the population (parent set) by applying the QEVSA $\frac{U}{2}$ times. |
| 5 | Crossover and mutation operations. |
| 6 | Produce $P_{st+1}$. |
| 7 | If the global result $F_{opt}$ is found, then halt, else go to 4. |

### III. Highly Constrained Quantum Genetic Algorithm

In this section, we extended the capability of the UQGA into the HCQGA, which can handle highly constrained optimization problems where no classical or quantum process can perform the search.

The HCQGA combines the strengths of the UQGA and penalty strategy to effectively determine the optimum extreme result of a highly constrained objective function, i.e., the algorithm specifically targets highly constrained optimization problems that are characterized by database beyond the computational power of both classical and quantum machines.

First, let's assume the task of optimizing the constrained objective function shown below,

$$maximize(G(\pmb{x})) \quad or \quad minimize(G(\pmb{x}))$$
$$s.t \quad e_y(\pmb{x}) \leq 0 \qquad \forall\, y = 1, \dots, w \quad , \qquad (1)$$
$$\pmb{x} \in X$$

The goal is to minimize (or maximize) the objective function $G(x)$ in relation to a number of inequality constraints $e_y(x)$, where $x$ is a $w$-dimensional vector with components $x_1$, $\ldots, x_i, \ldots, x_p$ that refers to the chromosome (possible candidate solution) of $G(x)$. The parameters $w$ and $p$ represent the total number of selected constraint parameters and the size of the variable $x$, respectively. Note that $x \in X$ describes the domain constraint for each variable $x_i$, i.e., the lower and upper bound of each variable $x_i$.

To solve the constrained optimization problem stated in (1), we transform it into an unconstrained one by applying the penalty strategy. Next, the process of solving (1) will be carried out in a similar manner as previously explained for UQGA.

The penalty method is a prevalent strategy utilized in genetic algorithms to address infeasible solutions to constrained optimization problems. Essentially, this method converts the constrained problem into an unconstrained one by penalizing infeasible solutions. This entails including a penalty term in the objective function for any constraint violation. For maximization problems, one writes the objective function with penalty term as,

$$eval(x) = G(x) + p(x), \qquad (2)$$

where $p(x)$ refers to the penalty function. To effectively solve maximization problems, certain criteria must be satisfied,

$$\begin{matrix} p(x) = 0 & if\ x\ is\ feasible \\ p(x) < 0 & otherwise \end{matrix}, \qquad (3)$$

In addition, the following should be satisfied,

$$|p(x)|_{max} \le |G(x)|_{min}, \qquad (4)$$

where, $|p(x)|_{max}$ and $|G(x)|_{min}$ denote the highest of $|p(x)|$ and lowest of $|G(x)|$ respectively, among solutions in the current population that are considered infeasible. For minimization problems, the penalty term should satisfy the following criteria,

$$\begin{matrix} p(x) = 0 & if\ x\ is\ feasible \\ p(x) > 0 & otherwise \end{matrix}. \qquad (5)$$

It is interesting to note that preparing the penalty term is not a straightforward task.

The HCQGA operates in a manner akin to the UQGA, albeit with a distinctive dissimilarity lying in the conversion of a constrained optimization problem into an unconstrained counterpart. An elucidating depiction of the operational process of the HCQGA can be observed in Figure 2, showcases a comprehensive flowchart.



Fig. 2. Working mechanism of the HCQGA.

The complexity of both strategies (HCQGA and UQGA) has been thoroughly examined in relation to number of generations $(G)$, population size $(U)$ and chromosome size $(C)$, providing a fair and comprehensive assessment. Therefore, the overall computational complexity of the algorithms is $O(GUC)$.

By replacing quantum initialization and selection stages instead of their classical counterparts, a significant reduction in the number of executed generations can be achieved. Furthermore, during the selection stage, quantum superposition enables the reduction of larger region sizes into smaller ones. Additionally, the size of the selected regions can impact the population size. Specifically, a larger region size can lead to a reduction in population size. It is worth noting that classically, the region size is denoted by $R$ and its processing time in classical computer is $O(R)$, whereas in quantum computing, the region size is represented as $log_2(R)$, and its processing time on a quantum computer is $log_2^3(R)$.

The chromosome size is contingent on the total size of the search space. Consequently, a larger search space results in a larger chromosome size. Moreover, storing the population size, requires $U$ classical register. In contrast, from the perspective of quantum computing, the power of quantum superposition enables the utilization of only one quantum register for storing the population size.

The running time of the HCQGA is contingent upon the execution time of the quantum computer. However, due to limited access to quantum computers, our simulations were conducted on classical computers. In the following, the

expression of the total running time of the HCQGA can be written as,

$$T_{HCQGA} = U.T_R + N_g(N).T_g, \qquad (6)$$

where $T_R$, $N_g$, and $T_g$ denote the processing time for the region $R$ by the quantum computer, the executed generation number _which depends strongly on the overall size of the database $N$_, and the running time required for one generation step, respectively.

## IV. MASSIVE MIMO SYSTEM

This section suggests an uplink multi-cell massive MIMO, in which the aim is to maximize the energy efficiency with respect to the desired bit rate of users. Followed by investigating how to implement the HCQGA in the proposed massive MIMO.

### A. Uplink Multi-Cell Massive MIMO Model

Let's consider an uplink scenario of a multi-cell system that implements massive MIMO technology. This system comprises $L$ cells sharing the same frequency band. In each cell, there exists a base station equipped with $M$ antennas, and $K_j$ single-antenna active users. It is important to note that active users have specific bit rate requirements, and we categorize users based on their target bit rates. The total number of target bit rate classes is denoted by $V$, with each class defined as $\eta^v$. Users are randomly distributed across the cells. Figure 3 illustrates the proposed uplink multi-cell massive MIMO system.



Fig. 3. A multi-cell massive MIMO system in the uplink configuration. Each cell is equipped with a single base station that has $M$ antennas. Within each cell, there are $K$ active users randomly distributed.
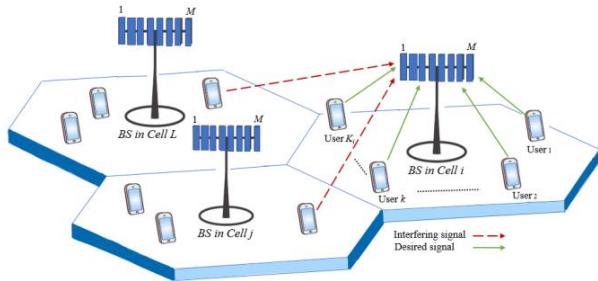
The channel gain, represented by $\beta_{ikm}$, characterizes both small and large-scale fading. Here, $i$, $k$, and $m$ correspond to the reference cell, the reference active user, and the received antenna at the base station in the reference cell $i$, respectively. Pilot contamination interference is not considered in this study, assuming that channel state information (CSI) is known in advance. The transmit power of the active user $k$ in the reference cell $i$ is denoted as $p_{ik}$.

The following equation represents the uplink transmission rate for a specific user in the system denoted by $Z$.

$$Z = \frac{\sum_m^M p_{ik}\beta_{ikm}}{\underbrace{\varepsilon \sum_{j \neq i}^L \sum_{l=1}^K \sum_{m=1}^M p_{jl}\beta_{jlm}}_{\xi_{inter}} + \underbrace{\varepsilon' \sum_{l \neq k}^K \sum_{m=1}^M p_{il}\beta_{ilm}}_{\xi_{intra}} + n_0}, \qquad (7)$$

where, the additive white Gaussian noise is denoted by $n_0$. Note that $\xi_{inter}$ describes the interference caused by users in other neighbouring cells, while $\xi_{intra}$ refers to the interference generated by users in the reference cell $i$. The scaling coefficients, $\varepsilon$ and $\varepsilon'$, are used to represent the interference ratios in the system, in which represent the interference ratio caused by interferer users in neighbouring cells (cells other than the reference cell $i$) and the interference ratio caused by interferer users within the same reference cell (cell $i$), respectively.

The expression of the achieved uplink transmission rate for the active user $k$ in the reference cell $i$ belonging to the target bit rate class $v$ as can be written as,

$$\eta_{ik}^v = W.log_2(1 + Z), \qquad (8)$$

where, $W$ represents the bandwidth utilized. Assuming all users share the same bandwidth $W$, we can determine the spectral efficiency of the reference cell $i$ as,

$$SE_i = \frac{\sum_{k=1}^K \eta_{ik}^v}{W}, \qquad (9)$$

where $\eta_i = \sum_{k=1}^K \eta_{ik}^v$ correspond to the transmission rate for all users in the reference cell $i$. For the purpose of this analysis, we assume that the bandwidth occupancy of all users is similar and, therefore, it is not considered. Assuming all cells have similar circuit power $P_i$. One expresses the energy efficiency of the reference cell $i$ as,

$$EE_i = \frac{SE_i}{\frac{1}{\gamma}\sum_k^K p_{ik} + L.P_i}, \qquad (10)$$

where $\gamma$ refers to the power amplifier efficiency coefficient. The goal of this research is to maximize the energy efficiency of cell $i$ while satisfying the desired bit rate for the active users. The optimization problem can be formulated as follows,

$$\begin{aligned} \max & \ EE_i \\ s.t \quad & \eta_{ik}^v \geq \eta^v \quad \forall i, k, v, \\ & p_{il} \geq 0 \ and \ p_{ik} \geq 0 \ \forall i, k, l \end{aligned} \qquad (11)$$

### B. How to apply the HCQGA in Massive MIMO

To solve (11), we employed the penalty approach defined in Section III. The new expression of the evaluate function can be written as,

$$eval(x) = EE_i + \varphi \sum_{k=1}^{K_i} max|0, R_{ik}^v - R^v|, \qquad (12)$$

where, $p(x) = \varphi \sum_{k=1}^{K_i} max|0, R_{ik}^v - R^v|$ refers to the global penalty term related to (12), while the parameter $\varphi$ denotes a penalty coefficient. The bottleneck exists in estimating properly the value of $\varphi$ in which depends on expression (12).

Several approaches are used to compute the penalty coefficient [37]–[40]. The choice of approach depends on the optimization problem and the required computational resources. Often, the fixed penalty coefficient method, which consists of setting up the penalty value manually and is easy to use and quick but lacks adaptability to data, is applied [37]. Another method is called the penalty parameter estimation strategy, which aims to estimate the penalty coefficient automatically based on the problem's characteristics. It is data-driven and improves generalization but can be computationally costly and assumption-dependent [39]. In addition to an adaptive penalty coefficient, which consists of modifying the penalty coefficient iteratively during the optimization process, it reduces the risk of regularization but requires complex implementation and tuning to avoid instability [40].

## V. SIMULATION RESULTS

In this section, we demonstrated the efficiency of the HCQGA in maximizing the energy efficiency of massive MIMO and reducing the computational complexity of the system through simulation evaluation. The classical constrained genetic algorithm (CCGA) was considered as a benchmark.

It is interesting to highlight that the overall number of active users and transmit power set _different possible power values that are considered as candidate values for transmitting a signal_ affect the overall search space, i.e., if the number of active users or the transmit power set increase, the database size also rises. For this sake, we constructed two simulation experiments, where each simulation analyzes the effect of each variant on the performance of the HCQGA and CCGA, i.e., the variant can be the overall number of active users or transmit power set.

### A. Simulation 1

To ensure the validity of our results, we implemented a well-established simulation framework. The simulation setup of the massive MIMO system consisted of one reference cell with a scaling factor of $\varepsilon' = 0.01$, surrounded by six neighboring cells. Each neighboring cell had an identical scaling factor of $\varepsilon = 0.00015$.

It is important to note that all cells, including the reference cell, contained only one base station with a fixed number of antennas, $M = 128$. The energy efficiency and the number of generations executed were computed for varying numbers of active users within the reference cell, ranging from 5 to 11 (denoted as $K$).

The study considers three distinct target bit rate classes, as indicated in Table 1. It should be noted that, for every given number of active users in the reference cell ($K_{ref}$), the total user population is randomly divided into two groups, each assigned to a specific target bit rate class. Furthermore, it is important to highlight that the computation of energy efficiency remains unaffected by the target bit rate of active users in the six neighboring cells.

In light of what has already been discussed, we concluded that converting the constrained objective function into an

TABLE I
SIMULATION PARAMETERS SETTINGS

| Symbol | Quantity | Value |
|--------|----------|-------|
| $B$ | bandwidth | 1 MHz |
| $\rho$ | static circuit power of BS | 1 W |
| $p_c$ | circuit power per antenna | 0.2 |
| $R^1$ | target bit rate class 1 | 100 Mbit/s |
| $R^2$ | target bit rate class 2 | 130 Mbit/s |
| $R^3$ | target bit rate class 3 | 150 Mbit/s |

unconstrained one and selecting the appropriate penalty coefficient depend tightly on the application, i.e., our suggested massive MIMO. In our investigation, we chose to select the optimal penalty coefficient manually. In order to approximate an appropriate setup for the penalty coefficient, we conducted a third simulation in parallel, using classical constrained exhaustive search (CCES) to identify the best optimal scenarios that maximizes the energy efficiency with the given target bit rate classes of the massive MIMO model. Since the CCES can handle small-scale databases, we employed a small number of users.

Figure 4 illustrates the energy efficiency utilized by the HCQGA and CCGA algorithms for different scenarios with varying numbers of total active users. As clearly shown in Figure 4, the energy efficiency usage with respect to the target bit rates is equal for both HCQGA and CCGA, across different numbers of active users. This demonstrates that both algorithms perform equally in maximizing energy efficiency while considering the target transmission rate of massive MIMO.



Fig. 4. Energy efficiency usage for different number of active users by the CCGA and HCQGA.

Figure 5 presents the average number of generations performed by the HCQGA and CCGA for different scenarios with varying numbers of total active users. In the comparison between the HCQGA and CCGA algorithms, it is noticeable that the CQGA exhibits a lower average number of generations performed compared to the CCGA. As a result, the computational complexity of the HCQGA is reduced in comparison to the

CCGA. Furthermore, as the overall number of active users (database size) increases, the HCQGA consistently maintains a lower computational complexity in comparison to the CCGA.



Fig. 5.  Average number of generations executed by the CCGA and HCQGA.

These findings highlight the computational advantages of the HCQGA algorithm in scenarios with a larger number of active users.

### B.  Simulation 2

The parameter setup of this simulation environment follows the previously described configuration in simulation 1. The objective of this simulation is to examine the influence of the transmit power set size on both the energy efficiency usage of the massive MIMO and the number of generations executed by the HCQGA and CCGA algorithms. The energy efficiency and the number of generations executed were computed for transmit power set size, ranging from 24 dBm to 42 dBm. The value of $M$ equals 128, while the value of $K$ equals 6.

Figure 6 shows the energy efficiency usage by the HCQGA and CCGA algorithms for different scenarios with varying transmit power set size. As seen from Figure 6, the energy efficiency usage with respect to the target bit rates, is found to be equal for both the HCQGA and CCGA algorithms across various transmit power set sizes. This observation affirms that both algorithms exhibit comparable performance in optimizing energy efficiency while accounting for the target transmission rate of the active users in massive MIMO.



Fig. 6.  Energy efficiency usage for different number of power set size for both the CCGA and HCQGA.

Figure 7 displays the average number of generations executed by the HCQGA and CCGA algorithms across different scenarios, considering various transmit power set sizes. In comparing the HCQGA and CCGA algorithms, it is evident that the HCQGA achieves a lower average number of generations performed compared to the CCGA, even when confronted with significant increases in transmit power set sizes. Consequently, the HCQGA exhibits a reduced computational complexity compared to the CCGA. Moreover, as the transmit power set sizes (database size) increases, the HCQGA consistently maintains a lower computational complexity relative to the CCGA.



Fig. 7.  Average number of generations executed by the CCGA and HCQGA.

### C.  Simulation 3

In the light of what have been discussed in the previous sections, it has been demonstrated that the quality of chromosomes during the initialization stage has a significant impact on the convergence speed of the genetic algorithm to the optimal result. This simulation investigates the effect of the coverage ratio $\frac{U.R}{N}$ on the runtime of the initialization stage of the HCQGA.

Figure 8 depicts the runtime required for various coverage values (%) during the quantum initialization stage of the HCQGA, when employing either classical or quantum computers. It is evident that as the coverage increases, the runtime of classical computers exhibits an exponential growth.



Fig. 8. The running time of the quantum and classical initialization stages.

In contrast, the quantum computer demonstrates significantly lower running times compared to the classical counterpart, while also maintaining consistently low runtimes despite increasing coverage.

Furthermore, Figure 9 illustrates the relationship between the region size and the average number of generations performed by the HCQGA. It is observed that as the region size increases, the number of executed generations drops.



Fig. 9.  The effect of increasing the region size on the executed average number of generations for the HCQGA.

## VI. Conclusion

This paper introduces the highly constrained quantum genetic algorithm (HCQGA) as a novel approach for finding the extreme value with respect to certain constraints in a vast and unsorted database, which surpasses the computational capacity of current available classical and quantum computers. To test the efficiency of this quantum strategy. We have investigated the maximization of the energy efficiency of an uplink multi-cell massive MIMO system with respect to the target bit rate of users. The HCQGA demonstrated accelerated convergence towards optimal maximum energy efficiency, outperforming the constrained classical genetic algorithm (CCGA) in terms of computational efficiency.

In future work, we are planning to extend the classical massive MIMO framework into a quantum counterpart, where entanglement-assisted quantum channels are considered for improving the performance of the overall bit rate and energy efficiency. Followed by demonstrating the theoretical results by building extensive simulations.

## References

[1] S. Imre and F. Balazs, *Quantum Computing and Communications: an engineering approach*. John Wiley & Sons, 2005.

[2] M. Mastriani, "Fourier Behind Entanglement: A Spectral Approach to the Quantum Internet," *Ann Phys*, vol. 534, no. 1, p. 2 100 296, Jan. 2022, **DOI**: 10.1002/andp.202100296.

[3] B. L. Márton, D. Istenes, and L. Bacsárdi, "Enhancing the operational efficiency of quantum random number generators," *Infocommunications Journal*, vol. 13, no. 2, pp. 10–18, 2021, **DOI**: 10.36244/ICJ.2021.2.2.

[4] L. Gyongyosi, L. Bacsardi, and S. Imre, "A Survey on Quantum Key Distribution," *Infocommunications Journal*, no. 2, pp. 14–21, 2019, **DOI**: 10.36244/ICJ.2019.2.2.

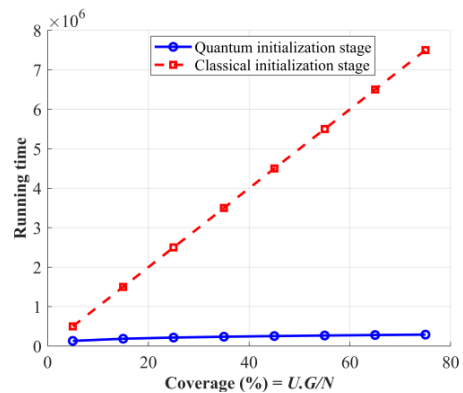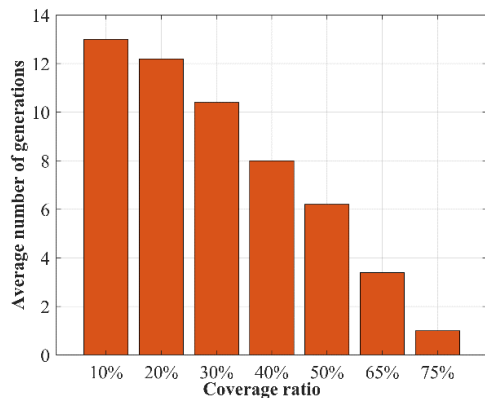[5] T. Bisztray and L. Bacsardi, "The Evolution of Free-Space Quantum Key Distribution," *Infocommunications Journal*, no. 1, pp. 22–30, 2018, **DOI**: 10.36244/ICJ.2018.1.4.

[6] D. Kobor and E. Udvary, "Optimisation of Optical Network for Continuous-Variable Quantum Key Distribution by Means of Simulation," *Infocommunications Journal*, vol. 12, no. 2, pp. 18–24, 2020, **DOI**: 10.36244/ICJ.2020.2.3.

[7] M. Czermann, P. Trócsányi, Z. Kis, B. Kovács, and L. Bacsárdi, "Demonstrating BB84 Quantum Key Distribution in the Physical Layer of an Optical Fiber Based System," *Infocommunications Journal*, vol. 13, no. 3, pp. 45–55, 2021, **DOI**: 10.36244/ICJ.2021.3.5.

[8] S. El Gaily and S. Imre, "Implementation of a constrained quantum optimisation method in resource distribution management with considering queueing scenarios," *International Journal of Communication Networks and Distributed Systems*, vol. 28, no. 2, p. 126, 2022, **DOI**: 10.1504/IJCNDS.2022.121194.

[9] G. Földes, "Techno-economic analysis on Mobile Network Sharing contribution to social welfare at 4G-5G area in Hungary," *Infocommunications Journal*, vol. 15, no. 1, pp. 87–97, 2023, **DOI**: 10.36244/ICJ.2023.1.9.

[10] A. Mihály and L. Bacsárdi, "Effects of selected noises on the quantum memory satellite based quantum repeaters," *Infocommunications Journal*, vol. 13, no. 2, pp. 19–24, 2021, **DOI**: 10.36244/ICJ.2021.2.3.

[11] M. Galambos and L. Bacsárdi, "Comparing Calculated and Measured Losses in a Satellite-Earth Quantum Channel," *Infocommunications Journal*, no. 3, pp. 14–19, 2018, **DOI**: 10.36244/ICJ.2018.3.3.

[12] M. Mastriani, "Fourier's Quantum Information Processing," *SN Comput Sci*, vol. 2, no. 2, p. 122, 2021, **DOI**: 10.36244/ICJ.2018.3.3

[13] A. Raja Basha, "A Review on Wireless Sensor Networks: Routing," *Wirel Pers Commun*, vol. 125, no. 1, pp. 897–937, Jul. 2022, **DOI**: 10.1007/s11277-022-09583-4.

[14] A. M. A. Sabaawi, M. R. Almasaoodi, S. El Gaily, and S. Imre, "New Constrained Quantum Optimization Algorithm for Power Allocation in MIMO," in *2022 45th International Conference on Telecommunications and Signal Processing (TSP)*, IEEE, 2022, pp. 146–149. **DOI**: 10.1109/TSP55681.2022.9851241.

[15] A. M. A. Sabaawi, M. R. Almasaoodi, S. El Gaily, and S. Imre, "MIMO System Based-Constrained Quantum optimization Solution," in *2022 13th International Symposium on Communication Systems, Networks and Digital Signal Processing (CSNDSP)*, IEEE, 2022, pp. 488–492. **DOI**: 10.1109/CSNDSP54353.2022.9907967.

[16] M. Almasaoodi, A. Sabaawi, S. El Gaily, and S. Imre, "Power Optimization of Massive MIMO Using Quantum Genetic Algorithm," in *1st Workshop on Intelligent Infocommunication Networks, Systems and Services (WI2NS2)*, Budapest University of Technology and Economics, 2023, pp. 89–94. **DOI**: 10.3311/WINS2023-016.

[17] S. el Gaily and S. Imre, "Quantum Optimization in Large Resource Management Systems," in *2019 IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, IEEE, Jul. 2019, pp. 1–5. **DOI**: 10.1109/SPAWC.2019.8815470.

[18] S. El Gaily and S. Imre, "Constrained Quantum Optimization for Resource Distribution Management," *International Journal of Advanced Computer Science and Applications*, vol. 12, no. 8, 2021. **DOI**: 10.14569/IJACSA.2021.0120806.

[19] M. Almasaoodi, A. Sabaawi, S. El Gaily, and S. Imre, "New Quantum Strategy for MIMO System Optimization," in *Proceedings of the 19th International Conference on Wireless Networks and Mobile Systems*, SCITEPRESS - Science and Technology Publications, 2022, pp. 61–68. DOI: 10.5220/0011305100003286.

[20] Y. Li, J. Ma, Z. Xie, Z. Hu, X. Shen, and K. Zhang, "A Scheduling Method for Heterogeneous Signal Processing Platforms Based on Quantum Genetic Algorithm," *Applied Sciences*, vol. 13, no. 7, p. 4428, Mar. DOI: 10.3390/app13074428.

[21] A. SaiToh, R. Rahimi, and M. Nakahara, "A quantum genetic algorithm with quantum crossover and mutation operations," *Quantum Inf Process*, vol. 13, no. 3, pp. 737–755, Mar. 2014, DOI: 10.1007/s11128-013-0686-6.

[22] J. S. Kim and C. W. Ahn, "Quantum strategy of population initialization in genetic algorithm," in *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, New York, NY, USA: ACM, Jul. 2022, pp. 439–442. DOI: 10.1145/3520304.3529010.

[23] U. Roy, S. Roy, and S. Nayek, "Optimization with quantum genetic algorithm," *Int J Comput Appl*, vol. 102, no. 16, pp. 1–7, 2014. DOI: 10.5120/17896-8732.

[24] J. Baghous, "5G system throughput performance evaluation using Massive-MIMO technology with Cluster Delay Line channel model and non-line of sight scenarios," *Infocommunications Journal*, vol. 13, no. 2, pp. 40–45, 2021. DOI: 10.36244/ICJ.2021.2.6.

[25] A. Popławski and S. Szott, "Using Dynamic Programming to Optimize Cellular Networks Modeled as Graphical Games," *Infocommunications Journal*, vol. 14, no. 4, pp. 62–69, 2022. DOI: 10.36244/ICJ.2022.4.9.

[26] X. Yu, Y. Du, X. Dang, S.-H. Leung, and H. Wang, "Power allocation schemes for uplink massive MIMO system in the presence of imperfect CSI," *IEEE Transactions on Signal Processing*, vol. 68, pp. 5968–5982, 2020. DOI: 10.1109/TSP.2020.3029404.

[27] T. Yu, S. Zhang, X. Chen, and X. Wang, "A Novel Energy Efficiency Metric for Next-Generation Green Wireless Communication Network Design," *IEEE Internet Things J*, vol. 10, no. 2, pp. 1746– 1760, Jan. 2023, DOI: 10.1109/JIOT.2022.3210166.

[28] M. R. Almasaoodi, A. M. A. Sabaawi, S. El Gaily, and S. Imre, "New Quantum Genetic Algorithm Based on Constrained Quantum Optimization," *Karbala International Journal of Modern Science*, vol. 9, no. 4, Oct. 2023, DOI: 10.33640/2405-609X.3325.

[29] A. Salh, N. S. M. Shah, L. Audah, Q. Abdullah, W. A. Jabbar, and M. Mohamad, "Energy-efficient power allocation and joint user association in multiuser-downlink massive MIMO system," *IEEE Access*, vol. 8, pp. 1314–1326, 2019. DOI: 10.1109/ACCESS.2019.2958640.

[30] J. Zhang, H. Deng, Y. Li, Z. Zhu, G. Liu, and H. Liu, "Energy Efficiency Optimization of Massive MIMO System with Uplink Multi-Cell Based on Imperfect CSI with Power Control," *Symmetry (Basel)*, vol. 14, no. 4, p. 780, Apr. 2022, DOI: 10.3390/sym14040780.

[31] H. Li, Z. Wang, and H. Wang, "Power allocation for an energy-efficient massive MIMO system with imperfect CSI," *IEEE Transactions on Green Communications and Networking*, vol. 4, no. 1, pp. 46–56, 2019. DOI: 10.1109/TGCN.2019.2948643.

[32] S. Imre, "Quantum existence testing and its application for finding extreme values in unsorted databases," *IEEE Transactions on Computers*, vol. 56, no. 5, pp. 706–710, 2007. DOI: 10.1109/TC.2007.1032.

[33] S. El Gaily and S. Imre, "Quantum Resource Distribution Management in Multi-task Environment," in *2019 14th International Conference on Advanced Technologies, Systems and Services in Telecommunications (TELSIKS)*, IEEE, Oct. 2019, pp. 364–367. DOI: 10.1109/TELSIKS46999.2019.9002340.

[34] S. Imre, "Extreme value searching in unsorted databases based on quantum computing," *International Journal of Quantum Information*, vol. 3, no. 01, pp. 171–176, 2005. DOI: 10.1142/S0219749905000700.

[35] S. El Gaily and S. Imre, "Quantum optimization of resource distribution management for multi-task, multi-subtasks," *Infocommunications Journal*, vol. 11, no. 4, pp. 47–53, 2019. DOI: 10.36244/ICJ.2019.4.7.

[36] J. F. Fitzsimons, "Private quantum computation: an introduction to blind quantum computing and related protocols," *npj Quantum Inf*, vol. 3, no. 1, p. 23, Jun. 2017, DOI: 10.1038/s41534-017-0025-3.

[37] M. S. Bazaraa, H. D. Sherali, and C. M. Shetty, *Nonlinear programming: theory and algorithms*. John Wiley & Sons, 2013.

[38] H. J. C. Barbosa and A. C. C. Lemonge, "An adaptive penalty scheme in genetic algorithms for constrained optimization problems," in *Proceedings of the 4th Annual Conference on Genetic and Evolutionary Computation*, 2002, pp. 287–294.

[39] S. N. Wood, "Modelling and smoothing parameter estimation with multiple quadratic penalties," *J R Stat Soc Series B Stat Methodol*, vol. 62, no. 2, pp. 413–428, 2000. DOI: 10.1111/1467-9868.00240

[40] B. Tessema and G. G. Yen, "A self adaptive penalty function based algorithm for constrained optimization," in *2006 IEEE international conference on evolutionary computation*, IEEE, 2006, pp. 246–253.

**Abdulbasit M. A. Sabaawi** has completed his bachelor's degree in computer and information Engineering at University of Mosul, Iraq in 2009, followed by a master's degree in communications and signal processing at Newcastle University, UK in 2014. He is currently a Ph.D. student at Budapest University of Technology and Economics (BME). His research focuses on OFDM, signal processing, Massive MIMO, quantum computing and communication, and quantum computing-based algorithms.

**Mohammed R. Almasaoodi**, a Ph.D. student at the Department of Networked Systems and services at the Budapest University of Technology and Economics (BME). he obtained his master's degree in computer engineering at University of Pecs in 2019 in Pecs, Hungary. he obtained his bachelor's degree Information engineering at Alnahrain university in 2009 in Baghdad, Iraq. His research interests are quantum computing, quantum communication, quantum computing-based algorithms, and telecommunication systems.

**Sara El Gaily**, a research fellow at the Department of Networked Systems and services at the Budapest University of Technology and Economics (BME). She received her Ph.D. in electrical engineering sciences in 2022 at Budapest University of Technology and Economics. She obtained her master's degree in electrical systems and renewable energies at Hassan II University in 2017 in Casablanca, Morocco. She obtained her bachelor's degree in computer and industrial electronics in 2014. Her research interests are quantum computing and telecommunication systems.

**Sándor Imre** [M'93] professor and Head of the Department of Networked Systems and services at the Budapest University of Technology and Economics (BME). He obtained dr. univ. degree in probability theory and statistics in 1996, a Ph.D. degree in 1999 and DSc degree from the Hungarian Academy of Sciences in 2007. He was elected the corresponding member of HAS in 2019. He acts as supervisor in the High-Speed Networks Laboratory since 1999. He is a member of the Doctoral Council of HAS. He was invited to join the eMobile Innovation Center of BME as an R&D director in 2005. His research interests include mobile and wireless systems, quantum computing and communications. Especially he has contributions to different wireless access technologies, mobility protocols and their game-theoretical approaches, reconfigurable systems, quantum computing-based algorithms, and protocols.

# Blockchain-Based, Confidentiality-Preserving Orchestration of Collaborative Workflows

Balázs Ádám Toldi, and Imre Kocsis

*Abstract*—Business process collaboration between independent parties is challenging when participants do not completely trust each other. Tracking actions and enforcing the activity authorizations of participants via blockchain-hosted smart contracts is an emerging solution to this lack of trust, with most state-of-the-art approaches generating the orchestrating smart contract logic from Business Process Model and Notation (BPMN) models. However, compared to centralized business process orchestration services, smart contract state typically leaks potentially sensitive information about the state of the collaboration, limiting the applicability of decentralized process orchestration. This paper presents a novel, collaboration confidentiality-preserving approach where the process orchestrator smart contract only stores encrypted and hashed process states and validates participant actions against a BPMN model using zero-knowledge proofs. We cover a subset of BPMN, which is sufficient from the practical point of view, support message-passing between participants, and provide an open-source, end-to-end prototype implementation that automatically generates the key software artifacts.

*Index Terms*—blockchain, BPMN, orchestration, collaboration, confidentiality, zero-knowledge proofs

## I. INTRODUCTION

In modern business science, *Business Process Management* (BPM) as a discipline [1] advocates process-focused thinking about internal activities and external collaborations to improve key performance indicators. Automating the execution of business processes is a key proposition of BPM and has been supported for a long time by various technical solutions [2]. Today, most of these, typically centralized, tools and services use the leading business process modeling standard, Business Process Model and Notation (BPMN) 2.0 [3] as a process definition language [4].

Distributed ledger technology (DLT), generally implemented on a blockchain basis, is widely recognized as a compelling platform to support the cross-organisational execution of business processes – even when the organisations cannot agree on a trusted (third) party as a middleman [5]. Blockchain-deployed smart contracts can impartially enforce the agreed-on sequences of activities and track sent and received messages. Smart contracts can also host data objects acted on by a process directly or anchor their changes in the blockchain via cryptographic commitments.

However, blockchain-assisted BPM is still a relatively new discipline – importantly, known BPMN-based solutions are inadequate from the privacy and confidentiality point of view. This paper presents a novel, collaboration confidentiality-preserving approach and end-to-end prototype tooling for the on-chain process orchestration of cross-organizational, BPMN-based collaborations using zero-knowledge proofs (ZKPs)[1]. Specifically, for a sufficient subset of BPMN, we present a transformation of the admissible state updates of BPMN process instances to programs of the ZoKrates [6] toolkit. We assemble state update validity provers from these programs for the participants and proof-verifying orchestrator smart contracts. We define an on-chain process state commitment update protocol, describe our open-source end-to-end implementation prototype[2] and evaluate practical viability.

Our contribution is novel from two aspects. First, to our knowledge, the confidentiality challenges of decentralized BPMN orchestration have not been addressed systematically and constructively yet. Second, we express BPMN execution as an incremental computation in a form amenable to commit-and-prove style zero-knowledge validation in smart contracts. This paves the way for further research on the computational representation of orchestrated BPMN execution against the continuously appearing ZKP advancements.

## II. MOTIVATION AND PROBLEM STATEMENT

BPMN is a standardized approach to visually and precisely express *how* business processes should be performed. BPMN is used in many domains – including finance, banking, manufacturing, healthcare, logistics and telecommunications – for capturing processes with well-defined sequences of regularly repeated activities. The BPMN standard defines several model types, *process*, *collaboration* and *choreography* being the most widely used ones. *Process (flow)* models are the simplest: these express the sequence, preconditions and exception handling of a single process performed by a single organization. Collaborations model the individual processes performed by collaborating parties – usually business entities – and their messaging-based interactions. Choreography diagrams focus solely on the message exchanges between collaborating entities.

Dept. of Measurement and Information Systems Budapest University of Technology and Economics Budapest, Hungary
(E-mail: balazs.toldi@edu.bme.hu, kocsis.imre@vik.bme.hu)

[1] This paper is based on the Scientific Student Association report submitted by Balázs Ádám Toldi to the 2022 competition at the Budapest University of Technology and Economics: https://tdk.bme.hu/VIK/sw8/Kollaborativ-munkafolyamatok-titkossagmegorzo
[2] Available at https://github.com/ftsrg/zkWF

## A. Decentralized orchestration

For over a decade, software tools have been available to assist with *process execution*. The more sophisticated ones track and *orchestrate* activities according to a BPMN model, register activity-related data and perform decision-making on further process evolution. However, centralized orchestration introduces a trusted orchestrator party requirement when we move beyond single-entity processes. With the emergence of blockchain and distributed ledger technology, the potential of decentralizing various aspects of cross-organizational collaboration has been recognized quite early.

Consider the BPMN car leasing collaboration model in Figure 1[3], where the internal processes of a car dealership, a leasing company and a financing bank must be coordinated to accept a leasing application. Individual executions of models are called *instances* and have an *instance state*. The coloring in Figure 1 demonstrates a state: green denotes that the dealership has completed insurance processing and the leasing company and the bank would be able to begin processing. However, for the leasing company to proceed, active downpayment checking (orange) must be finished, then the downpayment filed, and a downpayment notification sent and received.

Using blockchain-deployed smart contracts that track *collaboration state*, the *execution enablement* and *execution obligation* of the activities of the parties can be enforced without a dedicated, trusted party. The transaction journal nature of blockchains can also ensure that the full trace is also stored in an immutable and irrepudiable way. While tracking the internal state of participant-internal processes on-chain is not always desirable, it is a valuable *option*; e.g., when decisions have to be made in a way verifiable by the other collaborating parties.

Orchestrating and journaling messages and collaborative data handling are two further collaboration aspects which can be improved with "blockchainification". In both cases, the orchestrator smart contracts usually only manage cryptographic (hash) commitments to externally handled messages and data modifications, to avoid storing sizeable data on-chain.

Tools and approaches exist to create orchestrator smart contracts from BPMN models (see Section III). However, no systematic solution exists to protect sensitive collaboration state information in the smart contract state from parties who can read the blockchain but do not participate in the collaboration. In our example, a leasing company may wish that its competitors do not see how many open cases they have, how long it takes to perform key steps in the process, or what lease rates they apply.

Fulfilling such requirements is a confidentiality challenge that contradicts core blockchain design principles. Blockchain nodes must be able to *validate* and *execute* incoming transaction requests to reach consensus on ledger updates, be those changes of the balances of a natively tracked cryptocurrency or state changes of deployed smart contracts. If the transaction details are made "incomprehensible" to the nodes, e.g., by off-chain encryption, they can't validate the preconditions for performing the transaction and compute state updates. For smart contracts, the dominant *cryptographic* answers to this dilemma are validating transactions with ZKPs and confidentiality-preserving execution using homomorphic encryption, with the prior being significantly better established currently.

## B. Problem statement: BPMN collaboration confidentiality

We set up our problem statement through a basic system model and the enumeration of required security properties. We target a simple form of collaboration confidentiality (see the properties below) under the assumption that it is not in the interest of any process participant to leak information about process instances; participants neither directly leak information nor help external parties to compromise confidentiality. This is one of the realistic models for our setting, even though the participants do not completely trust the *actions* of each other. We will touch briefly on stronger models in Section IX.

*1) Basic assumptions and terminology: participants* wish to collaborate in the execution of an *instance* of a previously agreed-on BPMN collaboration definition. All other parties are *process external*. All participants have a cryptographic key pair for signature-based authentication and process activity authorisation. The underlying process model is public knowledge, but the public keys are shared only between the participants. We assume the absence of private key compromises.

For the underlying blockchain, we assume complete integrity (no successful attack on the consensus) and, for the sake of simplicity, deterministic finality (accepted blocks do not get retracted). Note that even blockchains with probabilistic block finality are usually quasi-deterministically final already at the time scale of a few blocks. On the other hand, process external parties have complete visibility of blockchain transactions. We treat the blockchain as *fair* – any transaction submitted by a participant is included in a block in a reasonable time, irrespective of concurrent transaction request load. While, in practice, blockchain platforms have strongly varying fault and threat models and sensitivity (see, e.g., [7]), these are basic assumptions of normal operational conditions. As a part of platform selection, security and dependability analysis should evaluate the risk of these assumptions not being met.

*2) System model:* the classic Business Process Orchestrator (BPO) middleware pattern [8] facilitates business process execution by providing a message broker and extending it with state management and persistent state storage. The solutions in the state-of-the-art closely match this pattern. (Technically, message passing is only *coordinated* and journaled by the smart contract their core.) The smart contract as a Process Controller [8] also performs authentication and authorization based on the BPMN model to ensure that the stored state sequence never deviates from the model semantics. We also aim to employ a blockchain-deployed smart contract as a BPO.

*3) Security properties:* we target a set of integrity, availability and confidentiality guarantees. Integrity and availability properties are already covered by the prior art; our contributions lie in establishing collaboration confidentiality, as defined by properties **C1** and **C2**, despite using smart contracts. BPO-SC refers to a per-process instance BPO smart contract.

---

[3] The model was created in the "Digitisation, artificial intelligence and data age workgroup" of the ongoing BME-MNB cooperation project. (MNB is the Central Bank of Hungary.). For legibility, the process in Figure 1 is slightly simplified; the whole model is available in our project repository.
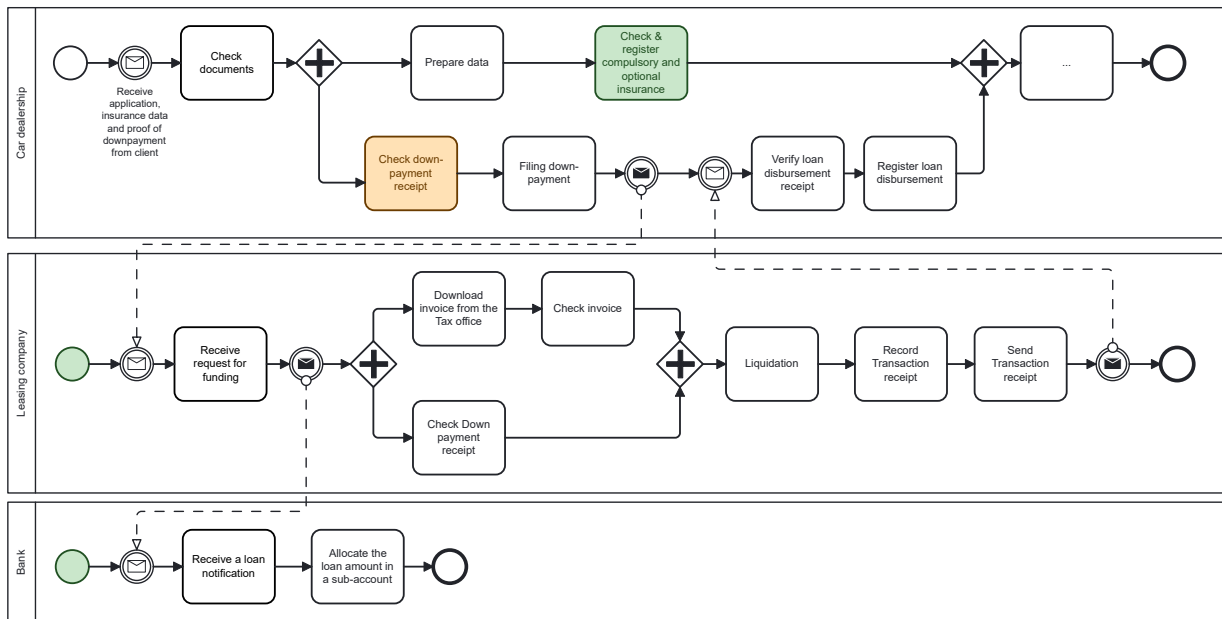
Fig. 1. Car leasing BPMN collaboration example (simplified for presentation). In the depicted example state, green denotes "completed"; orange is "active".

**I1**: The state traces enforced by the BPO-SC always adhere to the operational semantics of the underlying BPMN model.

**I2**: Process-external parties cannot influence the BPO-SC state.

**I3**: Process state updates can be initiated only by the participants authorized by the model and instance state.

**A1**: No external party can influence authorised participants' ability to perform state updates in a bounded time.

**A2**: No participant can influence the ability of authorized participants to perform state updates in a bounded time.

**A3**: Participants can always learn the trace and current state.

**C1**: External parties cannot determine participant identities.

**C2**: No external party can learn more about the trajectory, timings and stepwise properties (e.g., process variables and message contents) of the trace during and after execution than the fact that an instance has been started.

## III. RELATED WORK

Smart contracts, as a rule, cannot be altered after deployment; thus, to minimize the probability of software faults, domain-specific languages and Model-Driven Engineering (MDE) are steadily gaining ground in smart contract development [9]. In our context, the established approach is a BPMN model to serve as a *specification*, and orchestrator smart contract logic is generated automatically from the model.

### A. Decentralized business process orchestration

Caterpillar [10] was the first open-source BPMN-to-Solidity compiler (Solidity is the primary smart contract development language for the Ethereum platform). Since its initial release, several forks have emerged. Some of these also come with an extended feature set, like Blockchain Studio [11], which

adds role management, or [12], which adds time constraints. Lorikeet [13] is a model-driven engineering approach that integrates assets into business processes. Lorikeet extends the BPMN 2.0 specification with support for asset registries and also transforms models into Solidity smart contracts. The smart contracts handle the orchestration of the process as well as interactions with the tokens. Chorchain [14] takes a BPMN *choreography* and generates an Ethereum smart contract that can be used to execute the model. ChorChain also includes a dedicated modeling tool. The same authors released two further tools: Multi-Chain [15] and FlexChain [16]. Multi-chain is similar to Chorchain, but it also supports Hyper-ledger Fabric [17]. FlexChain can only produce Solidity smart contracts, but the user can also define a ruleset for each choreography. If a condition in the ruleset is met, then an off-chain processor will perform its underlying action.

Our analysis showed that the process state and trace are easily recoverable from the process manager smart contracts for *all* the tools above.

### B. Commit-and-prove ZKP with smart contracts

Zero-Knowledge Proofs (ZKPs) are cryptographic methods to prove the validity of various statements without revealing any additional information [18]. ZKP verification in a smart contract requires a scheme with "single-shot" message passing from prover to verifier; in this work, we rely on zk-SNARKs, a family of *noninteractive*, and also *succinct* (small and cheaply verifiable proofs) ZKPs. We use the ZoKrates toolkit as a ZKP front-end with a high-level programming language [6]. ZoKrates currently supports the Groth16 [19], GM17 [20] and Marlin [21] proving schemes.

Our contribution implements a commit-and-prove approach. In commit-and-prove schemes, a party first commits to an
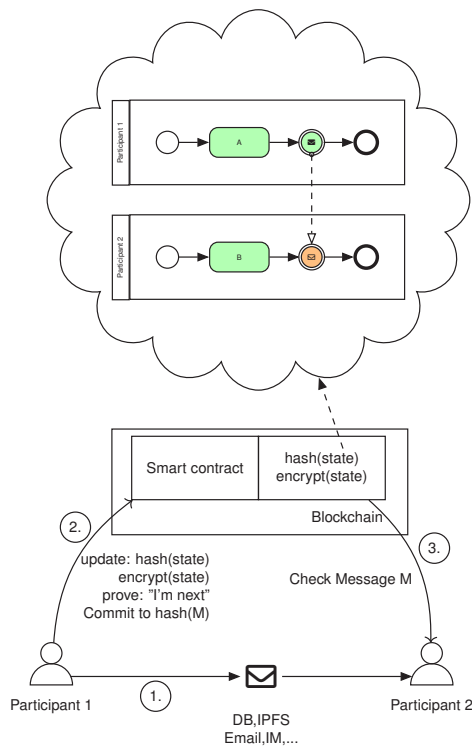
Fig. 2. Overview of the zkWF protocol

input and, possibly later, proves some predicate about the input – without revealing it [22]. This is a widely used pattern in the smart contract-based application of zero-knowledge proofs. Recent surveys on ZKP schemes, technologies and applications can be found in [23] and [24].

## IV. A CONFIDENTIALITY-PRESERVING APPROACH

The fundamental difference of our approach from [10–16] is that instead of storing the process state on-chain in an easily interpretable form in an orchestrator smart contract, we store encrypted states and cryptographic state commitments and accept update proposals on the presentation of ZKPs over the current and proposed commitment. The approach relies on two key conceptual components: our zkWF ("zero knowledge WorkFlow") protocol and what we call "zkWF programs".

### A. The zero-knowledge WorkFlow (zkWF) protocol

The zkWF protocol is a hash commitment style protocol that allows the participants of a business process to follow and step the execution of a business process. Figure 2 presents a high-level overview.

At the centre of the scheme is a smart contract instance on a blockchain for an instance of a collaboration model. This smart contract stores and manages the state of the collaboration – as specified by the underlying BPMN model – in an encrypted and a hashed form.

During process execution, the collaborating parties can send messages to each other by off-chain means ①. These are

captured in the underlying process specification as intermediate message throw and capture events; our state commitment scheme includes commitments to the message hashes.

When a participant wishes to update the state stored in the smart contract – that is, to "step the process" –, it has to create a ZKP that the proposed state transition is valid. This new state includes the hash of the message they sent beforehand if the step involves message sending. It sends the new state hash commitment, the encrypted new state and the ZKP proof of state transition validity to the smart contract as a blockchain transaction ②; the smart contract updates its state only if it can successfully check the ZKP.

When the execution arrives at a point where a participant receives a message in the next stage of the execution, the receiving party checks the hash and only accepts (and proceeds with its part of the collaboration) if the hashes match ③.

Participant authentication is tied to proving private key ownership in the ZKPs. The public keys are defined over the participant group-shared process model as a parameterization. These are cooperation-private, "application-level" key pairs; on pseudonymizing platforms, such as Ethereum, updater identity can and should be masked by using independent, single-use transaction source addresses (i.e., public keys).

Additionally, we require the participants to have a common means for encrypting and decrypting stored state ciphertexts. The protocol does not constrain the encryption used.

The protocol can be realized straightforwardly on a wide range of DLTs; we provide an implementation for Ethereum and Hyperledger Fabric [17]. While the updates and the contract state are unintelligible to parties outside the collaboration, statistical and model trace analyses of the update sequences are still a threat. We enable mitigations by including a "fake" update transaction variant (no actual state update), which all participants are authorized to use.

### B. zkWF programs

zkWF programs are generated from BPMN specifications and serve as a bridge between process definition and proof computation/verification. A zkWF program is a ZoKrates program that, for a given BPMN model instance (parameterized model), can decide whether a given actor is authorized to execute a state transition in a given execution state. We use the zkWF program to generate the zero-knowledge proofs and proof verification code for the orchestrator smart contract.

### C. Workflow and toolchain

We created an end-to-end toolchain prototype for our approach, as depicted in Figure 3.

In the *modeling phase*, a BPMN model is annotated with metadata for process instantiation, and our interpreter-translator creates the corresponding zkWF program.

In the *synthesis phase*, the ZoKrates toolkit is used to set up the *prover key* and *verifier key* and generates the verifier smart contract in Solidity. We created novel support for generating verifier code for Hyperledger Fabric in Java. We also created the code generation facilities for both platforms' state commitment management part of the smart contracts.
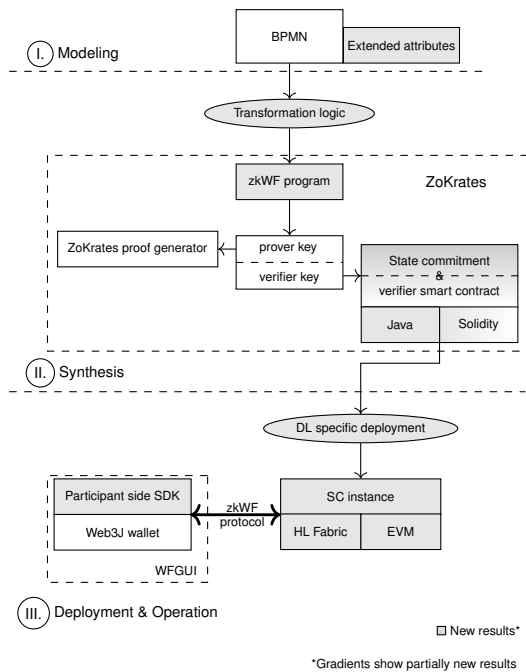
Fig. 3. Toolchain overview

Some secret values used when creating zk-SNARK prover and verifier keys are considered "toxic waste": an adversary can use them to break the scheme, e.g., forge fake proofs. Thus, security relies on the waste having been deleted. The associated risk can be mitigated by using a reliable party for the key generation or performing so-called multi-party trusted setup ceremonies, where a (large) group of actors assembles the keys. In this case, security requires only at least one of them to delete the waste. Such ceremonies tend to be complicated and thus can pose a problem for by-program setup. Universal schemes also exist (e.g., [21]), where the results of a single program-agnostic ceremony can be used to derive program-specific keys publicly and securely. Choosing the right approach requires deployment-specific risk analysis; ZoKrates supports all of the above.

For the *deployment phase*, we created automation facilities for deployment to Ethereum (and other blockchains using a compatible RPC API); and an SDK and GUI for the client side. Here, we integrate the ZoKrates toolkit as a proof generator.

## V. BPMN SUBSET AND EXECUTION SEMANTICS

This paper targets the Basic Modeling Elements of BPMN 2.0 [3, p. 28], the core subset of the specification, with the restrictions that regarding events, we interpret only message throw and catch ones (among participants) and do not support sub-processes and data objects. We argue that this element set is already sufficient for practical applications. Statistical evidence [25] shows that the usage frequency of the 50 constructs in the BPMN specification follows a Zipfian distribution; we cover elements used at least in $\sim 25\%$ of the models in [25]. This is also the empirically established "Common Core of

BPMN" in [25] with the addition of messaging between participants. Our earlier example showcases the supported element set (except for "exclusive gateways" for process variable-based choice paths and "lanes" for further subdividing pools).

### A. BPMN extensions and structural constraints

We introduce two extended attributes for BPMN elements. `zkp:publicKey` separates the tasks of different participants by attaching a participant-specific public key to a pool, a lane, or a *participant-executable element* (activities, message throws and catches). `zkp:variables` applies to *activities* and declares process instance global variables, and that that activity may write the variable (reads are allowed for all activities). These variables can be used in boolean expressions for exclusive gateways.

Some constraints apply to the structure of the BPMN models, which are currently admissible in our scheme.

- Gateways must be binary (two incoming/outgoing edges).
- Activities are *atomic*; i.e., subprocesses are not supported.
- The model must be acyclic (no loops).

We plan to eliminate these constraints in the future; the required modifications of the state representation and the zkWF program construction are largely incremental.

### B. State representation

Our notion of process instance execution state encompasses the following aspects (for the specific encoding in zkWF programs, please refer to the report and the implementation).

- A vector $v$ of the current state of executable elements
- The current values of *global variables*
- Hashes of the messages already sent in the process

Let $M = (V, E, T)$ be a process model, where $V$ is the set of non-flow model elements, $E$ is the set of model edges (flows), and $T \subset V$ is the set of all executable elements in the business process. Then, $v$ is a vector of $|T|$ size and $\forall v_i \in v$ can have one of the following three values:

- 0 (Inactive) – The element has not been reached yet
- 1 (Active) – The element is ready to be executed or is being executed by a participant
- 2 (Completed) – The execution of the element has been completed

This state set is a subset of those in the standard activity lifecycle [3, p. 428] and serves as a reasonable simplification, as the main focus of the work described here is exploring the confidential execution aspect. Note that correctly implementing the full lifecycle is a significant software engineering effort, even in the centralized setting. Also, BPMN users tend to apply a similar simplified view during modeling, as the more sophisticated state aspects require experience and limit the ease of model understanding.

### C. Capturing token passing semantics

BPMN 2.0 models have straightforward, token flow-based standard execution semantics: start events create tokens that move around as execution progresses. Parallel gateways split

and join tokens. To support a different ZKP use case, [26] introduces a technique for representing valid BPMN execution state changes by enumerating the possible composite token marking deltas of the elements upon stepping the process. Specifically, [26] introduces an array $P$, where each element of $P$ is a list of token change and element identifier pairs. We construct a similar $P$ array under the token passing semantics and embed it into the zkWF program to enable checking whether a proposed state update is valid from the BPMN execution logic point of view. Our $P$ array to describe one-step token marking changes for a model $M$ consists of 3-tuples with elements from the set $\mathcal{N}$:

$$\mathcal{N} = (+1, -1) \times T) \cup \{(0, -1)\} \qquad (1)$$

For $T$, we apply a simple integer encoding; the $-1$ in the "no-token-change" pair second set is a don't care placeholder. Especially under our binary gateway condition, which is currently necessary to ensure reasonable proof computation times, it is straightforward to enumerate the admissible changes based on the BPMN model. For example, let's consider activities $a, b, c \in T$. $a$ continues in a parallel gateway, which proceeds to $b$ and $c$. When $a$ transitions from "Active" to "Completed" and $b$ and $c$ from "Inactive" to "Active", the following token marking change happens: $((-1, a), (+1, b), (+1, c)) \in \mathcal{N}$. The complete logic can be found in the referenced report.

## VI. ZKWF PROGRAM AND PROTOCOL DESIGN

A zkWF program is a ZoKrates program shared among the participants, with which process participants prove that a business process state transition they propose is allowed. In ZKP terms, the participants are the *provers*, and the orchestrator smart contract is the *verifier*.

ZoKrates programs have public as well as private inputs, and an output. Private inputs are only visible to the prover; public inputs are visible to the prover and the verifier, and they are necessary to verify proofs. In our case, the current commitment and the proposed one act as public inputs. Private inputs are more varied; only some are shared across the participants (e.g., the cleartext of the current state).

The key current deficiency of our scheme is that our proofs do not include showing the congruence of the on-chain stored state ciphertexts and the public state (hash) commitments. Combining established encryption algorithms with zk-SNARKs is hard; advances are being made (see, e.g., [27]), but these haven't appeared in any of the leading zk-SNARK frameworks yet as vetted and reusable "gadgets".

We apply the following measures to this deficiency. An additional part of our public input (and blockchain-stored data) will be a *signature commitment*: the current hash commitment and the *previous* hash commitment signed by the last acting party (using their application-level cryptographic identity). Should a participant erroneously or maliciously commit a ciphertext that does not hash to the stated, proven and accepted commitment, this signature ensures that the offending participant can be irrepudiably identified by the other collaborating parties.

Although several partially mitigative and corrective schemes can be built on this measure, we introduce the weakening



Fig. 4. The basic computation model of zkWF programs

assumption that the irrepudiable identifiability of participants halting execution this way is a sufficient disincentive.

### A. zkWF computation model

Figure 4 illustrates the structure of the generated zkWF programs. For hashing, we use SHA-256; application-level signing uses the EdDSA implementation from the ZoKrates standard library (both widely used, NIST-standard algorithms). The *private* inputs of zkWF programs are as follows.

- $s_{current}$ - the current state of the process (subsec. V-B)
- $r_{current}$ - random salt for hashing $s_{current}$ (32 bits)
- $s_{new}$ - the updated ("stepped") process state
- $r_{new}$ - new randomness, for hashing $s_{new}$
- $pk$ - public EdDSA key of the participant (subsec. IV-A)
- $sk$ - private EdDSA key of the participant

The *public* inputs ($\|$ denotes concatenation):

- $h_{current} = hash(s_{current} \| r_{current})$
- $S_{new} = sig(h_{current} \| h_{new})$

$sig$ denotes signing by the party proposing the new hash commitment in the concatenation. Given these inputs, the following steps are performed.

1) Checking the group-shared secret current state and randomness against the public hash commitment to ensure ongoing integrity.
2) Checking that no illegal state transition is being proposed through $s_{new}$ at the *process logic* level.
3) Checking the new signature commitment given as a public input (based on $pk$ and $sk$) and checking the authorization of the participant for the business process step.
4) The program outputs the hash of the new state.

Most aspects of the computational model are straightforward; we only expand on the important details of BPMN model encoding and the state change validity checking logic.

### B. BPMN model encoding and state change validation

The BPMN model logic is carried over into the zkWF program by a precomputed $P$ array (Section V-C). To check whether the correct paths are proposed for exclusive gateways,

the expressions on the sequence flows after the gateways are also encoded in the program as assertions. Message passing and variable write permission checks are addressed similarly.

Regarding the executable element state vector, the program compares $v_{current}$ and $v_{new}$ from $s_{current}$ and $s_{new}$. If the two are the same, the "change" is accepted (as a "step" under our fake update mechanism). Four or more differences (pairwise comparisons at the same indices) in the vectors are considered invalid. Otherwise, we construct a $3 \times 3$ matrix $A$ with the initial value

$$A = \begin{bmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{bmatrix} \quad (2)$$

Then, for the $j$-th difference ($j \in 0 \ldots 2$) at position $i \in 0 \ldots |T| - 1$ in the vectors, we apply the following updates:

- $v_{current}[i] = 1$ & $v_{new}[i] = 2 \Rightarrow A[j] \leftarrow [-1, i]$
- $v_{current}[i] = 0$ & $v_{new}[i] = 1 \Rightarrow A[j] \leftarrow [1, i]$
- $v_{current}[i] = 0$ & $v_{new}[i] = 2 \Rightarrow A[j] \leftarrow [1, i]$

Any other combination of $v_{current}$ and $v_{new}$ values is invalid. If $P$ contains an element with the rows of $A$, then token passing-wise, the proposed state change is valid, as we essentially decoded the activity token marking changes ($\pm 1$) from the activity state changes: 0 - Inactive $\rightarrow$ 1 - Active $\rightarrow$ 2 - Completed.

Parallel gateway ends ("joins") induce an additional check: a transition from a state where *not both activities before the gateway are completed* to one where *both are* also requires that the activity after the gateway gets activated. State change validation also includes checking write permissions for global variables and contrasts the evaluation of arithmetic expressions with the proposed path for exclusive gateways.

Finally, the message-handling validation logic involves two major validation aspects. First, a message hash must be provided when a participant wants to mark a Message Throw event as "completed". We assume the actual message to be passed off-chain. Second, when a participant wants to mark a Message Catch event as "completed", we must ensure that the corresponding Message Throw event is also marked as completed. The receiver contrasts the message with the hash value; if this fails, we assume that the further steps are either captured in the process logic or the sender and receiver coordinate corrective transmission off-chain.

### C. The zkWF protocol

The protocol flows through the orchestrator smart contract and is simple in light of the earlier sections. The smart contract state contains the following elements:

- $h_{current} = hash(s_{current}||r_{current})$
- $C_{curr}^{enc} = enc(s_{current}, r_{current})$
- $S_{current} = sig(h_{prev}||h_{current})$

where $enc$ denotes encryption with the group encryption key and method (see Section IV). Update request transactions of the smart contract carry the following arguments:

- $h_{new} = hash(s_{new}||r_{new})$
- $C_{new}^{enc} = enc(s_{new}, r_{new})$

- $S_{new} = sig(h_{current}||h_{new})$
- $p(h_{current}, S_{new}, h_{new})$

The last argument is a ZKP of the correspondence of $h_{current}$, $S_{new}$ and $h_{new}$, under the shared zkWF program. The orchestrator smart contract checks the validity of this proof before accepting the smart contract state change carried by the other arguments.

### D. Side-channel attack protections

Public BPMN models facilitate side-channel attacks on confidentiality. Our work until now aimed to ensure that the trace steps of the BPMN finite automaton remain unintelligible to the external observer; however, the number and timings of transitions still carry information. Most BPMN models are simple enough to infer a usable probability distribution of possible states and traces from just these observations.

Constant-time execution and delay randomization are two apparent protection options, though both introduce artificial delays. Consider a constant-time token passing ring schedule with dummy operations as our already established scheme. For $n$ participants, we determine a suitable time quantum $t$ with which it is acceptable to wait for $(n-1)t$ to delay the "posting" of any state change. During process execution, at the beginning of the $i$-th epoch, participant $i \bmod n$ checks whether it needs to send a state update transaction. If yes, it does; if not, it issues a "fake update" transaction. After terminating the process, a long fake update stream is advisable. As long as enough participants meet their fake update obligations and adhere to their epochs, external observers only see a heartbeat-like stream of uninterpretable transactions and can determine even the time of termination only with low probability.

### VII. Security properties

The presented approach addresses the security requirements defined in Section II-B as discussed in this section.

### A. Integrity

Property I1 holds in the sense that we carefully implement a strict subset of BPMN semantics, but we acknowledge that future work should create an explicit proof of conformance. I2 holds due to application-level cryptographic authentication; I3 due to cryptographic authentication and the very simple sub-logic of enabling activities and message operations.

### B. Availability

A1 holds due to I2 and the blockchain fairness assumption – which is mild for high-throughput public and cross-organizational blockchains. A2 holds only under the disincentive assumption of Section VI. However, the assumption is not strong for domains with a credible threat of legal or regulatory action (e.g., finance). A participant can also perform a denial of service attack with a constant stream of malicious fake updates. The disincentive assumption applies here, too, but fake update regimen-dependent defences can also be introduced in the smart contract (e.g., epoch schedule enforcement). A3 holds due to a smart contract accounting for state and trace and the blockchain platform assumptions.

### C. Confidentiality

The C1 guarantee has two layers. At the platform level, all transactions can originate from single-use addresses on pseudonymizing platforms – e.g., Ethereum. In Hyperledger Fabric, the Identity Mixer protocol suite for transactor anonymization and unlinkability can be used similarly. At the application level, transaction payloads and smart contract states contain only hashed, signed and encrypted data. Hashing is straightforward; for the signed content, note that EdDSA signatures do not provide a way to recover the signer's public key from the signature or to determine whether the same key was used to sign two different messages. For the encrypted state, if not a single, group-shared secret is used, an application should choose an encryption scheme where the participant keys cannot be recovered.

C2 depends on external data and transaction uninterpretability, which flows from the cryptographic measures, and transaction unlinkability, which also relies on the measures for C1. It also requires sufficient side-channel protection, for which we have at least one strong (not necessarily efficient) option.

## VIII. IMPLEMENTATION, TESTING AND PERFORMANCE

The ZoKrates toolkit is a central component in our framework; the current implementation uses version 0.7.13[4]. ZoKrates was the ZKP toolkit with the best-fitting programming language and ZKP scheme support during our research.

### A. Code generation

Our code generator, implementing the transformation logic denoted in Figure 4, is a custom development in Kotlin. This component generates a zkWF program from an XML-serialized BPMN model, relying on ZoKrates template files. First, the model is encoded, as we outlined earlier; then, it generates the code for the described stages of computation and checks.We also generate the orchestrator smart contracts for EVM-based blockchains (Solidity version 0.8.0) and Hyperledger Fabric (Java "chaincode").

### B. Client side

We created a simple participant-side SDK, which wraps ZoKrates and incorporates the Web3J wallet library. We also created a TornadoFX-based desktop GUI application ("Work-Flow GUI") for testing and demonstration purposes. The GUI supports all key participant-side actions: monitoring a process manager smart contract for changes, retrieving state, creating process step proposals, computing their witnesses and proofs, and submitting update proposals.

WFGUI also incorporates a process modeller for our BPMN subset and extensions through an embedding of bpmn-js[5]; supports testing through preassembled smart contract call sequences; and supports process manager smart contract deployment to Ethereum-based blockchains. A demonstrational video is available in our repository.

---

[4] See https://github.com/Zokrates/ZoKrates/releases
[5] See https://bpmn.io/toolkit/bpmn-js/

### C. Functional testing

We assembled a suite of *simple test cases*, based on the test model suites of the tools cited in Section III. BPMN model size and complexity influence zkWF program size and complexity, which, in turn, determine proof computation times and on-chain verification costs. To evaluate the practical feasibility of our approach, the leasing model from Section II was used as our *representative test case*.

### D. Performance evaluation

In addition to functional testing (compliance with model semantics, proper enforcement of authorization aspects and proper handling of compliant/noncompliant proofs), we used our test suite to evaluate key performance metrics of the approach. Performance tests were performed on a desktop PC (AMD Ryzen 7 2700, 16 GB of DDR4 memory).

In Ethereum, smart contract execution steps, measured in "gas", incur a cryptocurrency cost, paid by the transaction-requesting user. For measurements of gas used, we used a private, one-node Ethereum test network with version 1.10.25 of geth, the official Go implementation of the Ethereum protocol. Blockchain-side efficiency measurements are largely irrelevant for Hyperledger Fabric, which has no "gas" notion and where the smart contract execution layer is highly resource-scalable. Table I summarizes the relevant size metrics of our test cases.

TABLE I
BPMN MODEL AND TEST CASE CHARACTERISTICS

| Case | Vertices | Edges | Executable | Size of $P$ | Scenarios |
|---|---|---|---|---|---|
| Test 1 | 5 | 4 | 3 | 3 | 3 |
| Test 2 | 9 | 10 | 5 | 7 | 9 |
| Test 3 | 8 | 8 | 4 | 4 | 4 |
| Test 4 | 6 | 5 | 2 | 3 | 2 |
| Test 5 | 14 | 12 | 10 | 10 | 10 |
| Repr. | 68 | 69 | 50 | 54 | 52 |

Table II summarizes the runtimes of the off-chain computations. Compilation and zk-SNARK setup were executed once; proving time is the sum of computing the witness and generating the proof, and we give an average over the scenarios. The measurements indicate that our approach is practically feasible for real-life models.

TABLE II
OFF-CHAIN COMPUTATION RUNTIMES

| Case | Compilation time | Setup time | Proving time avg. |
|---|---|---|---|
| Test 1 | 27.22 s | 129.58 s | 55.0 s |
| Test 2 | 48.32 s | 182.80 s | 88.67 s |
| Test 3 | 28.55 s | 129.69 s | 53.40 s |
| Test 4 | 27.14 s | 128.82 s | 53.21 s |
| Test 5 | 30.74 s | 133.44 s | 54.10 s |
| Repr. | 81.02 s | 187.33 s | 122.47s |

Table III summarizes the gas costs of smart contract deployment and smart contract calls in the zkWF protocol. Note that although the representative model is 5-6 times larger than the simple ones, the smart contract call gas cost is only moderately higher. As the hashes, signatures, and proofs have a fixed length, gas usage variability is driven by the size of the

encrypted version of the current state. In the measurements, we use state cleartext instead of ciphertext to eliminate the impact of the not-constrained encryption.

TABLE III
GAS USAGE ON ETHEREUM

| Case | Deployment gas usage | Update gas usage avg. |
|------|---------------------|----------------------|
| Test 1 | 2,098,786 gas | 490,507 gas |
| Test 2 | 2,098,990 gas | 497,780 gas |
| Test 3 | 2,098,498 gas | 493,705 gas |
| Test 4 | 2,078,071 gas | 503,817 gas |
| Test 5 | 2,161,039 gas | 491,783 gas |
| Repr. | 2,408,635 gas | 548,898 gas |

Due to the novelty of our approach, it is comparable with the state of the art only in gas costs. Deployment is on par with, or is better than, the existing solutions. However, the cost of updating the state is significantly higher; ChorChain uses about 92,905 gas on average for each message and Caterpillar is similar to ChorChain.

This "confidentiality premium" is certainly not acceptable on the Ethereum mainnet. Still, it can be argued that the high gas price on the mainnet has "priced out" all use cases that were not strictly crypto-financial years ago. On the other hand, at the time of this writing, on multiple alternative EVM-based public blockchains, the gas costs of our operations translate to fractions of 1 USD. Additionally, our approach has evident usage potential on purpose-created, permissioned, cross-organizational blockchains; in this case, the gas cost is a technical consideration and low enough to allow for dozens of transactions per block under customary block gas targets. Lastly, we store encrypted state on-chain "only" to fulfil requirement A3 the simplest way; highly available off-chain data storage with blockchain-based integrity assurance is a common technique.

## IX. THREATS TO VALIDITY AND FUTURE WORK

We see compliance with BPMN operational semantics as a non-negligible threat to validity, especially after our planned future extension of the supported BPMN subset. For the approach presented in this paper, we only tested compliant behavior and not formally prove it; this remains future work.

Impractical proof time for much larger BPMN models is also a threat. We plan to introduce the capability to handle *hierarchical* process models. We expect that we can instantiate orchestrator smart contracts for sub-processes in a way that coordinates the commitment-management across the levels, but controls proof obligation complexity by requiring proof generation only for a limited-size model part for each update.

While the ring schedule "fake updates" approach is evidently correct for adhering participants (and, we surmise, for mostly adhering participants), side-channel protections is an open line of research. We plan to analyse the ring schedule scheme under various participant failure models and compare it with delay randomization schemes. Metrics for measuring the guaranteed level of protection through fake updates are necessary, too. Differential privacy metrics worked out for publicly observable messaging settings with a "hide-in-the-crowd" approach similar to ours [28] promise to be adaptable.

Lastly, we note that there are stronger versions of our collaboration confidentiality model through additional *inter-collaborator confidentiality constraints*; it is an interesting question how our approach can be extended to such settings.

## X. CONCLUSION

In this paper, we presented a collaboration confidentiality-preserving approach for the smart contract-based orchestration of business collaborations, captured as BPMN 2.0 models. Our protocol is a novel, and to our knowledge, first-of-its-kind solution, which we validated functionally as well as evaluated from the resource usage and gas cost points of view. We also described a full toolchain prototype which we made available as open-source software.

## ACKNOWLEDGMENT

## REFERENCES

[1] W. M. Van Der Aalst, M. La Rosa, and F. M. Santoro, "Business process management: Don't forget to improve the process!" *Business & Information Systems Engineering*, vol. 58, no. 1, pp. 1–6, 2016. **DOI**: 10.1007/s12599-015-0409-x

[2] S. Pourmirza, S. Peters, R. Dijkman, and P. Grefen, "A systematic literature review on the architecture of business process management systems," *Information Systems*, vol. 66, pp. 43–58, 2017. **DOI**: 10.1016/j.is.2017.01.007

[3] Object Management Group, "Business Process Model and Notation (BPMN), Version 2.0." [Online]. Available: https://www.omg.org/spec/BPMN/2.0/

[4] M. Chinosi and A. Trombetta, "BPMN: An introduction to the standard," *Computer Standards & Interfaces*, vol. 34, no. 1, pp. 124–134, 2012. **DOI**: 10.1016/j.csi.2011.06.002

[5] J. Mendling et al., "Blockchains for Business Process Management - Challenges and Opportunities," *ACM Trans. Manage. Inf. Syst.*, vol. 9, no. 1, 2018. **DOI**: 10.1145/3183367

[6] J. Eberhardt and S. Tai, "ZoKrates – Scalable Privacy – Preserving Off-Chain Computations," in *2018 IEEE International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)*, 2018, pp. 1084–1091. **DOI**: 10.1109/Cybermatics 2018.2018.00199.

[7] C. Cachin and M. Vukolić, "Blockchain consensus protocols in the wild," in *31st International Symposium on Distributed Computing (DISC 2017)*, 2017. [Online]. Available: https://drops.dagstuhl.de/opus/volltexte/2017/8016/pdf/LIPIcs-DISC-2017-1.pdf

[8] I. Gorton, *Essential Software Architecture*, 2nd ed. Springer Berlin, Heidelberg, 2011.

[9] Y. Ait Hsain, N. Laaz, and S. Mbarki, "Ethereum's Smart Contracts Construction and Development using Model Driven Engineering Technologies: a Review," *Procedia Computer Science*, vol. 184, pp. 785–790, 2021. **DOI**: 10.1016/j.procs.2021.03.097

[10] O. López-Pintado, L. García-Bañuelos, M. Dumas, I. Weber, and A. Ponomarev, "Caterpillar: A business process execution engine on the ethereum blockchain," *Software: Practice and Experience*, vol. 49, no. 7, pp. 1162–1193, 2019. **DOI**: 10.1002/spe.2702

[11] L. Mercenne, K.-L. Brousmiche, and E. B. Hamida, "Blockchain Studio: A Role-Based Business Workflows Management System," in *2018 IEEE 9th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON)*, 2018. pp. 1215–1220 **DOI**: 10.1109/IEMCON.2018.8614879.

[12] A. Abid, S. Cheikhrouhou, and M. Jmaiel, "Modelling and Executing Time-Aware Processes in Trustless Blockchain Environment," in *Risks and Security of Internet and Systems*, ser. LNCS, 2020. pp. 325–341, **DOI**: 10.1007/978-3-030-41568-6_21.

[13] Q. Lu et al., "Integrated model-driven engineering of blockchain applications for business processes and asset management," *Software: Practice and Experience*, vol. 51, no. 5, pp. 1059–1079, 2021. **DOI**: 10.1002/spe.2931

[14] F. Corradini et al., "ChorChain: A Model-Driven Framework for Choreography-Based Systems Using Blockchain," in *Proc. of the 1st Italian Forum on Business Process Management (ITBPM)*, 2021, pp. 26–32.

[15] ——, "Model-driven engineering for multi-party business processes on multiple blockchains," *Blockchain: Research and Applications*, vol. 2, no. 3, p. 100 018, 2021. **DOI**: 10.1016/j.bcra.2021.100018

[16] ——, "Flexible Execution of Multi-Party Business Processes on Blockchain," in *2022 IEEE/ACM 5th International Workshop on Emerging Trends in Software Engineering for Blockchain (WETSEB)*, 2022, pp. 25–32. **DOI**: 10.1145/3528226.3528369.

[17] E. Androulaki et al., "Hyperledger Fabric: a distributed operating system for permissioned blockchains," in *Proceedings of the Thirteenth EuroSys Conference*, 2018, pp. 1–15. **DOI**: 10.1145/3190508.3190538.

[18] ZKProof Community, "ZKProof Community Reference," 2022. [Online]. Available: https://docs.zkproof.org/reference.pdf

[19] J. Groth, "On the Size of Pairing-based Non-interactive Arguments," in *Advances in Cryptology – EUROCRYPT 2016: 35th Annual International Conference on the Theory and Applications of Cryptographic Techniques*, 2016, pp. 305–326. **DOI**: 10.1007/978-3-662-49896-5_11.

[20] J. Groth and M. Maller, "Snarky signatures: Mini- mal signatures of knowledge from simulation-extractable SNARKs," in *Advances in Cryptology – CRYPTO 2017: 37th Annual International Cryptology Conference*, 2017, pp. 581–612. **DOI**: 10.1007/978-3-319-63715-0_20.

[21] A. Chiesa, Y. Hu, M. Maller, P. Mishra, P. Vesely, and N. Ward, "Marlin: Preprocessing zkSNARKs with Universal and Updatable SRS," Cryptology ePrint Archive, Paper 2019/1047, 2019. [Online]. Available: https://eprint.iacr.org/2019/1047

[22] D. Bennaroch, M. Campanelli, D. Fiore, J. Kim, J. Lee, H. Oh, and A. Querol, "Proposal: Commit-and-Prove Zero-Knowledge Proof Systems and Extensions," https://docs.zkproof.org/standards/proposals, presented at the 4th workshop of the ZKProof Community, 19-29 April 2021, online.

[23] X. Sun, F. R. Yu, P. Zhang, Z. Sun, W. Xie, and X. Peng, "A Survey on Zero-Knowledge Proof in Blockchain," *IEEE Network*, vol. 35, no. 4, pp. 198–205, 2021. **DOI**: 10.1109/MNET.011.2000473

[24] J. Partala, T. H. Nguyen, and S. Pirttikangas, "Non-interactive zero-knowledge for blockchain: A survey," *IEEE Access*, vol. 8, pp. 227 945–227 961, 2020. **DOI**: 10.1109/ACCESS.2020.3046025

[25] M. z. Muehlen and J. Recker, "How Much Language Is Enough? Theoretical and Practical Use of the Business Process Modeling Notation," in *Seminal Contributions to Information Systems Engineering: 25 Years of CAiSE*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 429–443. **DOI**: 10.1007/978-3-642-36926-1_35.

[26] T. Aivo, "Zero-Knowledge Proofs for Business Processes," Master's thesis, Univ. of Tartu, 2020.

[27] J. Lee, J. Choi, J. Kim, and H. Oh, "SAVER: SNARK- friendly, Additively-homomorphic, and Verifiable Encryption and decryption with Rerandomization," Cryptology ePrint Archive, Paper 2019/1270, 2019. [Online]. Available: https://eprint.iacr.org/2019/1270

[28] I. A. Seres, B. Pejó, and P. Burcsi, "The Effect of False Positives: Why Fuzzy Message Detection Leads to Fuzzy Privacy Guarantees?" in *Financial Cryptography and Data Security*, ser. LNCS. Cham: Springer International Publishing, 2022, pp. 123–148. **DOI**: 10.1007/978-3-031-18283-9_7.

**Balázs Ádám Toldi** received his BSc in computer engineering in 2023 from the Budapest University of Technology and Economics (BME). Currently, he is an MSc student at BME, with a primary specialization in cybersecurity and a secondary specialization in critical systems.

**Imre Kocsis** received his PhD from the Budapest University of Technology and Economics (BME) in 2019. Currently, he serves as a senior lecturer and leading blockchain researcher at the Critical Systems research group of the Dept. of Measurement and Information Systems of BME. He leads the activities of the group in conjunction with the Hyperledger Foundation and the university's participation in the European Blockchain Services Infrastructure (EBSI) network.

# Balancing Information Preservation and Data Volume Reduction: Adaptive Flow Aggregation in Flow Metering Systems

Adrian Pekar, Laszlo A. Makara, Winston K. G. Seah, and Oscar Mauricio Caicedo Rendon

*Abstract*—The critical role of network traffic measurement and analysis extends across a range of network operations, ensuring quality of service, security, and efficient resource management. Despite the ubiquity of flow-level measurement, the escalating size of flow entries presents significant scalability issues. This study explores the implications of adaptive gradual flow aggre- gation, a solution devised to mitigate these challenges, on flow information distortion. The investigation maintains flow records in buffers of varying aggregation levels, iteratively adjusted based on the changing traffic load mirrored in CPU and memory utilization. Findings underscore the efficiency of adaptive gradual flow aggregation, particularly when applied to a specific buffer, yielding an optimal balance between information preservation and memory utilization. The paper highlights the particular significance of this approach in Internet of Things (IoT) and contrasted environments, characterized by stringent resource constraints. Consequently, it casts light on the imperative for adaptability in flow aggregation methods, the impact of these techniques on information distortion, and their influence on network operations. This research offers a foundation for future studies targeting the development of more adaptive and effective flow measurement techniques in diverse and resource-limited network environments.

*Index Terms*—adaptive computing, gradual flow aggregation, flow measurement, data reduction, performance optimization

## I. INTRODUCTION

As the digital age progresses, networks become increasingly complex, connecting myriad devices and applications. Each of these applications and services possesses unique requirements, making the task of managing and understanding network traffic of paramount importance. This management hinges on effective *traffic measurement* and *analysis*. The significance of this task goes beyond mere data monitoring. It provides essential data for operations such as SLA compliance evaluation, QoS provisioning, intrusion detection, and traffic management [1].

The predominant approach to understanding network traffic is through *flow level* analysis [2]. Conceptually, a "flow" represents a series of packets, moving from a source to a destination, that share some common attributes. It is a foundational concept that serves as the basis for the more complex methodologies discussed in this paper. However, as networks expand in complexity and size, grappling with the *expanding magnitude of flow records* becomes an uphill task. Several strategies aim to manage this via *adaptive aggregation of flow records*, wherein flow records are dynamically consolidated based on certain parameters to streamline the data without losing significant detail.

However, as traffic volumes surge, particularly under heavy loads, the very technique of adaptive flow aggregation can alter the richness and reliability of traffic information. This distorted information can directly impede the efficiency of network operations, making it crucial to understand flow aggregation's true implications. Particularly, with the growing emphasis on the Internet of Things (IoT) and similar network environments where resources are often constrained, understanding this balance becomes critical. As edge computing starts to overshadow cloud-based solutions, especially in the IoT landscape, ensuring minimum resource consumption while maintaining data richness becomes indispensable.

Given this backdrop, our research delves into the realm of adaptive flow aggregation, probing its effects on flow size distortion—a critical metric determining the reliability of traffic management. Specifically, our exploration involves analyzing flow records organized across buffers with varying degrees of aggregation, which are dynamically generated based on traffic loads reflected through CPU and memory utilization.

One of our pivotal observations reveals that employing adaptive flow aggregation can achieve a desirable trade-off. Specifically, a certain degree of aggregation resulted in a marginal information loss of just $2.42\%$, a compromise that significantly optimizes resource utilization. This insight is particularly vital in settings like IoT where there is a pressing need to maximize data integrity while operating within constrained resources. The distinguishing contributions of this paper encompass:

(*i*) A deep dive into how differential aggregation levels, particularly in light of fluctuating traffic volumes, influence the fidelity of flow information.

(*ii*) The introduction of the SEE indicator, a novel metric to quantify information biases induced by flow aggregation. SEE provides a rigorous gauge of estimation errors in flow size—a cornerstone for gauging the repercussions of adaptive aggregation on flow data.

A. Pekar and L. A. Makara are with the Department of Networked Systems and Services, Faculty of Electrical Engineering and Informatics, Budapest University of Technology and Economics, Budapest, Hungary.

A. Pekar is also with the HUN-REN-BME Information Systems Research Group, Budapest, Hungary.

W. K. G. Seah is with the School of Engineering and Computer Science, Victoria University of Wellington, Wellington, New Zealand

O. M. Caicedo Rendon is with the Department of Telematics Engineering, Engineering Telematics Group, Universidad del Cauca, Popayán, Colombia

Corresponding author: A. Pekar (e-mail: apekar@hit.bme.hu)

(*iii*) The exploration into the aftereffects of adaptive flow aggregation on flow size distortion, enriched by SEE's integration.

Our research accentuates the pressing need for truly adaptive flow aggregation methodologies, underscoring their implications on data fidelity and overall network operations. We aspire that our insights catalyze the development of refined, adaptive flow measurement techniques tailored for diverse and resource-challenged network settings.

The rest of this paper is structured as follows: Section II provides a succinct background on flow measurement, aggregation, and its adaptive variants. Section III reviews relevant literature in the domain. In Section IV, we outline our research design, highlighting our novel gradual aggregation scheme, adjustment strategies, preliminaries, and measurement techniques. Section V presents the findings of our study, analyzed to enhance comprehension of our approach. Discussions on these results are in Section VI, while Sections VII and VIII delve into their broader significance. Finally, Section IX offers reflections on our research.

## II. Background

This section provides a brief yet comprehensive overview of network traffic flow aggregation in a context relevant to the scope of this paper. From exploring the dynamics of network traffic flow, through presenting an overview of the process of flow aggregation, to considering the principle of adaptive flow aggregation, this section sets the stage for the subsequent methodology of our study.

### A. The Dynamics of Network Traffic Flow

Presently, the most prevalent method for network measurement is the collection of traffic data at the flow level, commonly known as flow export [3]. The term *flow* denotes a group of packets that possess a shared *key* and pass through a specific observation point within a determined period [4]. This shared *flow key* is typically characterized by a five-element tuple, including the source and destination IP addresses, source and destination ports, and the protocol.

Traffic data, such as *flow features*—for instance, all packets related to a specific flow quantified in bytes—and the corresponding flow keys, often designated as *flow properties*, are contained in the *flow records*. Network management tasks rely on the analysis of these flow records, where the most crucial data is procured from the packet headers encapsulated in the flow features. Analyzing these records includes descriptive calculations like determining the minimum, mean, standard deviation, and maximum of the flow sizes along with the packet inter-arrival time statistics.

In a conventional flow-based measurement scenario, the flow records traverse through several platform components positioned at different layers, as illustrated in Figure 1. The procedure typically involves organizing the captured packets into flows in a *flow cache*, after which they are intermittently exported to a data collector. The collector then either stores the data in a data store or forwards them directly for further analysis and visualization [2].
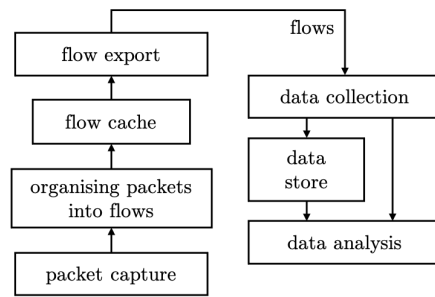


Fig. 1. General architecture of a flow measurement platform.

However, the enormous quantity of measurement data presents several challenges [5] as they traverse through the platform's components, right from packet capture to data analysis. Making sense of the intricate, ever-expanding data for operational purposes is a daunting task. The processing essential for both online and offline analysis poses considerable challenges, often due to a lack of sufficient resources to manage flow-related tasks efficiently, accurately, and sustainably [6].

### B. The Impact of Tuple Variations on Potential Flow Entries

In the realm of IP flow measurement, the flow cache plays a pivotal role by aggregating packets into flows using various combinations of a five-tuple. This aggregation results in flow entries, which provide a comprehensive view of network traffic. The nature and number of these entries can vary based on the specifics of the tuple in use and the practicalities of network operations.

The five-tuple used for flow measurement in IP networking typically consists of the following:

1) *Source IP Address*: IPv4 addresses have a 32-bit length, leading to $2^{32}$ possible addresses. IPv6 addresses, on the other hand, have a 128-bit length, leading to $2^{128}$ possible addresses.
2) *Destination IP Address*: Same as above.
3) *Source Port*: There are 65,536 possible port numbers (ranging from 0 to 65,535).
4) *Destination Port*: Same as the source port.
5) *Protocol*: The IP protocol field is 8 bits, so there are $2^8$ or 256 possible values. However, not all 256 values are used in practice. Common protocols include TCP, UDP, and ICMP. For simplicity, we use the full range of 256 possible values.

In Table I, the potential combinations of these elements showcase the vast number of unique flow entries one could encounter in IP networking. Evidently, the magnitude of possible combinations is colossal, especially when factoring in the expansive address space of IPv6. However, as more fields from the five-tuple are excluded, the number of combinations decreases substantially.

It is important to note that these figures represent a theoretical maximum. In practice, IP flow measurement on a specific network segment does not capture all these combinations because of the inherent nature of network design and purpose,

TABLE I
THE TOTAL NUMBER OF POSSIBLE TUPLES FOR DIFFERENT
COMBINATIONS

| Tuple Combination | IPv4 | IPv6 |
|---|---|---|
| Src IP, Dst IP, Src Port, Dst Port, Proto | $1.1 \times 10^{28}$ | $1.2 \times 10^{79}$ |
| Dst IP, Src Port, Dst Port, Proto | $1.1 \times 10^{21}$ | $1.5 \times 10^{43}$ |
| Src Port, Dst Port, Proto | $1.1 \times 10^{13}$ | $1.1 \times 10^{13}$ |
| Dst Port, Proto | $1.7 \times 10^{7}$ | $1.7 \times 10^{7}$ |
| Proto | 256 | 256 |

hence observing only a subset of potential traffic patterns. This highlights the nuanced and comprehensive scope of flow measurements in IP networks, drawing a distinction between theoretical possibilities and real-world observations.

### C. The Process of Flow Aggregation

Flow aggregation serves as a popular method to streamline the handling of measurement data. It primarily aims to combine several underlying measurement data points into a single, unified record. The merging process is governed by properties whose values may remain static or fluctuate over time. It should be noted that flow records are themselves representations of aggregated traffic, with flow features compiled from information sourced from packet headers.

However, flow aggregation stands as a higher abstraction level where the properties of multiple flow records are synthesized according to specific criteria, yielding a more concise representation of the original data. This aggregation results in consolidated flow records that maintain the overall traffic properties and characteristics, albeit with a broader granularity and diminished detail.

While aggregation invariably results in flow records with lesser information richness about network traffic, the information loss incurred is substantially less than that observed in sampling-based methods [7]. In certain traffic dynamics, these methods can overlook entire flows, typically those composed of a solitary packet or very few packets. This limitation has led to the development of existing flow aggregation techniques [8], [9] which employ partial aggregation of flow records deemed non-essential. The process is facilitated by *gradual flow key reduction*—that is, a step-by-step reduction in the count of flow properties that function as flow keys—in order to safeguard crucial information and preserve accuracy.

However, without the careful calibration of aggregation intensity, achieving the appropriate granularity of flow-level information becomes challenging. Ensuring optimal levels of aggregation becomes essential in maintaining the delicate balance between data manageability and the preservation of pertinent information.

### D. The Principle of Adaptive Flow Aggregation

Adaptive systems leverage iterative adjustments to their parameters based on predefined criteria to ensure that the system operates optimally or as close to optimal conditions as possible. Adaptive aggregation adheres to a similar principle—modifying aggregation according to one or more criteria to enhance system operation. Such criteria could include the flow size (measured in bytes), CPU or memory utilization by the flow measurement process, link bandwidth utilization of the capture device, or a combination thereof.

For instance, if flow aggregation is adjusted in response to CPU utilization, the number of processed packets can be perceived as directly proportional to the CPU utilization by the flow metering tool. When the traffic volume increases, the count of captured packets also escalates, subsequently intensifying the resource utilization of flow measurement. In contrast, a decline in traffic leads to a decreased packet count, thereby reducing the associated resource utilization. This reciprocal relationship allows the definition of different network traffic load levels—such as low, moderate, and high. These levels serve as markers for estimating how the traffic load is registered by flow measurement, enabling the adjustment of flow aggregation accordingly.

Adaptive flow aggregation proves beneficial as it applies a relatively less destructive degree of aggregation contingent upon the release of additional resource capacity. As a result, more informational value is preserved when fluctuations in resource capacity permit. This advantage is particularly notable in the context of constrained IoT operations, which often face limitations in processing capability and memory. Furthermore, as recent advancements in IoT solutions demonstrate a trend towards transitioning computational logic from the cloud to the edge [10], the need for low resource consumption while preserving information value is increasingly crucial.

Contrarily, traditional flow aggregation with static operation introduces a consistent level of distortion to flow-level information, irrespective of resource availability. This can result in unavoidable information loss during memory overflow events. Notably, although adaptive flow aggregation employs more aggressive merging under high loads, it helps to maintain information integrity even when the flow cache overflows (i.e., when the memory is entirely filled). Despite the fact that this approach might result in lower information granularity in the aggregated flow records, it ensures no information is lost in overflow situations. Furthermore, the process of gradual flow key reduction is specifically designed to mitigate the erosion of relevant information.

### III. RELATED WORK

The domain of network management has witnessed extensive research, but only a few studies have specifically honed in on flow aggregation. An earlier work, Aguri, offers a distinct approach to aggregation-centric traffic profiling suitable for real-time, long-term, and wide-area traffic monitoring [9], [11]. Unlike conventional methods that rely on predefined filter rules to classify traffic types, Aguri aggregates low-volume flows until they become distinctly identifiable. This ensures even minor traffic types are not overlooked. With its capacity to generate concise profiles spanning source and destination addresses and protocols, Aguri adeptly monitors traffic, spots anomalies, and counters threats like DDoS attacks.

Diving deeper into the nuances of flow aggregation, Cheng *et al.* [12] proposed an Aggregation Flow Measurement

(AFM). This scheme reconfigures traffic clusters leveraging the quintessential five fields of a fine-grained flow. Central to this method is dynamic sampling that adeptly recalibrates in response to traffic shifts. Instead of direct packet value recordings, it banks on estimates, optimizing CPU resource allocation. Adding finesse to this approach is a secondary process that prioritizes heavy-tailed flows during flow information updates. This dual-pronged strategy ensures rich information capture and impeccable estimation accuracy while judiciously pruning smaller flows to streamline flow cache.

In assessing these methodologies, it is clear that adaptation has been progressively factored into flow aggregation techniques. While Aguri offers a unique perspective on traffic profiling without considering real-time adaptation, AFM brings resource utilization, specifically CPU dynamics, into the fold. However, packet sampling entails a loss of granularity, potential inaccuracies, challenges in rare event detection, inconsistent accuracy across flows, and a lapse in bursty traffic information capture. The added layer of flow sampling in AFM aggravates these issues as flows, once organized, could be counterproductively discarded.

Gradual Flow Key Reduction by Irino et al. [8] presents a more refined approach. It orchestrates flows based on predetermined criteria—like the volume of transported octets in descending order. Only flows surpassing a user-defined significance threshold are retained; the rest undergo iterative aggregation. This proactive strategy curtails network congestion and safeguards upper-tier platform components, as visualized in Figure 1. However, despite its valuable insights, this method does not fully embrace 'adaptation.' It misses reacting to pivotal aspects such as resource utilization by measurement tools or the fluidity of network traffic. Undoubtedly, there is still a broader spectrum of criteria, like the full range of network traffic dynamics, that might further refine the granularity and efficiency of flow-level data processing.

In a notable shift, recent studies increasingly integrate flow aggregation with Heavy Hitter (HH) flow detection [13], [14]. Central to this detection is a threshold demarcating HHs from their non-HH counterparts. Recognizing the challenges posed by anomalous traffic patterns on measurement tools, Hu et al. [15], [16] introduced a dynamic strategy utilizing adaptive flow aggregation. At the heart of the methodology is the insight that a significant proportion of network attacks manifest as non-HH patterns, predominantly generating a plethora of short-lived flows. The primary application of this method finds relevance in scenarios like Denial of Service (DoS) attacks that target the same destination IP address, and worm attacks with identical source IP addresses. In response to these situations, their approach combines a two-dimensional hash table structure with a tiered clustering system. Once sorted, the traffic within these clusters is aggregated into metaflows, mitigating issues related to memory and export bandwidth, while preserving the accuracy of legitimate flows.

Pekar et al. [17] introduced a technique inspired by Gradual Flow Key Reduction but extended its application to the domain of HH detection. This method tailors the aggregation of flow records to match the characteristics of the traffic and the specific objectives of the monitoring process. For instance, this could include anomaly detection or flow-based accounting of transferred data volumes. In this scheme, HHs are maintained in the flow cache, while non-HHs are aggregated in dedicated buffers, following a hierarchy defined by flow key precedence. However, the work examines the adaptive aggregation of flow records exclusively through the prism of reduction efficacy in relation to HH detection, thereby also neglecting to explore the performance efficacy, particularly in relation to information distortion and resource consumption.

Building upon the theme of HH detection in software-defined data center networks, Bi et al. [18] sculpted a dynamic threshold, examining the nexus between elephant thresholds and network traffic dynamics in data centers. Drawing parallels with the optimal receive system of baseband signal transmission, they unearthed an equilibrium between positive and negative false rate detections by studying the overlap of two flow probability distribution curves.

In a similar vein, Wang et al. [19], [20] differentiated flows into short and long categories, applying a fluid threshold. The unique twist in their approach is the distinct management of these flows: short flows are shepherded by distributed algorithms, while long ones are entrusted to centralized solutions. This method is further enhanced by capitalizing on end-hosts for precise flow tagging and curtailing overhead through centralized algorithms.

Lastly, Liu et al. [21] leveraged the Dynamical Traffic Learning algorithm. This tool facilitates real-time dynamic configuration of threshold values, ensuring swift and efficient identification of HH flows with minimal latency and overhead.

Pivoting from flow record data to table rule management, Saha et al. [22] hones in on flow table rule aggregation. Focusing on QoS, flow table entry aggregation, and IoT, their adaptive scheme reduces flow-rules volume without sacrificing IoT traffic's QoS. A key-based mechanism empowers user choice over OpenFlow match-fields. Balancing QoS path selection and switch flow-table use, the "Best-fit" heuristic adaptively selects a QoS path to minimize network flow-rules.

Navigating the intricate orchestration of the small flows emblematic of modern mobile core networks and IoT, Minh et al. [23] introduces the 'flow tree', a controller-based binary search tree-like structure mirroring the OpenFlow switch table. By exploiting the wildcard of a flow's dstIP, the method manages the tree responsively to network shifts, resulting in a flow table that is more efficient and trimmer than standard OpenFlow tables.

Taking a step further, Phan et al. [24] developed a mechanism to optimize traffic flow monitoring in SDN-based networks. It adjusts flow table entries in SDN switches based on the detailed traffic information required by systems such as intrusion detection or traffic engineering. Instead of a fixed threshold, the method uses a machine-learning algorithm to determine optimal entry limits. This continual assessment ensures SDN switch performance remains optimal.

Lastly, in an endeavor to streamline table occupancy in SDNs, Jia et al. [25] showcases a flow-table aggregation

strategy. Through dynamic address and port rewriting, the method aggregates multiple same-destination flows from varied sources into a singular flow entry. This approach drastically reduces core-layer SDN switch table occupancy, proving effective even in environments with dispersed IP address allocations. The method can operate in both software-defined IPv4 and IPv6 networks, though it does not provide explicit adaptability features.

In reflection, while foundational groundwork in flow measurement has been laid by prior research, there is a distinct evolution from static criteria to more dynamic, HH detection-driven approaches. Nevertheless, predominant focus areas have been network security and measurement data reduction. The overarching challenge has been to comprehensively blend adaptability, precision in aggregation, and nuanced flow handling. The landscape, enriched by these advancements, still beckons for solutions that holistically address the multifaceted challenges of modern network environments.

Our study steps into this gap with an innovative approach. Through the adaptive multi-buffer flow measurement strategy, we aim to navigate beyond the traditional confines of prior research. Our methodology fuses adaptability with structured aggregation, ensuring that flow data is managed judiciously and resourcefully. By doing so, we intend to shed light on the impact of adaptive flow aggregation on information bias and flow measurement instrumentation, providing a roadmap for improved network traffic management. In essence, our approach endeavors to set a new course in extracting meaningful insights from flow data, even in environments marked by resource variability. To our knowledge, this paper is the first effort in bridging this particular domain gap.

## IV. METHODOLOGY

This section delves into our methodology designed to study adaptive gradual flow aggregation. We pay special attention to the aggregation and adaptation techniques, foundational assumptions, primary propositions, and our evaluative metrics.

### A. Gradual Flow Aggregation Scheme

Anchoring our strategy is the *Gradual Flow Key Reduction* mechanism [8]. This technique systematically and progressively aggregates flow data, trimming flow key elements as data navigate a sequence of buffers. It pivots on two fundamental components: the *flow key precedence* and the *multi-buffer structure*.

*1) Flow Key Precedence:* The linchpin of our strategy is the preordained *Flow Key Precedence*. This order determines the hierarchy of flow key element reduction. As we traverse each level of reduction, we witness flow records of differing granularity—starting from the most detailed, progressively becoming more aggregated.

*2) Multi-buffer Structure:* Our multi-buffer structure embodies a delicate equilibrium between data granularity and resource constraints, setting forth a well-considered data reduction route. Under the assumption of a standard five-tuple flow key precedence as 'protocol' > 'src_port' > 'dst_port'

> 'src_ip' > 'dst_ip', we detail the flow granularity across buffers as follows:

- *Main Flow Cache (Buffer B0)*: Retains the entire flow key set, guaranteeing maximum information fidelity.
- *Buffer B1*: Initiates aggregation by removing 'dst_ip', leaving behind 'protocol', 'src_port', 'dst_port', and 'src_ip'.
- *Buffer B2*: Proceeds further by excluding 'src_ip', encompassing only 'protocol', 'src_port', and 'dst_port'.
- *Buffer B3*: Excludes 'dst_port' next, safeguarding 'protocol' and 'src_port'.
- *Buffer B4*: At its apex of aggregation, only the 'protocol' remains.

Even though rooted in the predominant five-tuple, our methodology displays inherent adaptability. We ensure that the buffer configuration aligns with the diversity of flow key elements, thereby providing versatility for various flow definitions.

*3) Gradual Flow Aggregation Operation:* The aggregation cycle commences by targeting the lowest-ranking flow key. Resultant flow records, housed in buffer $B_1$, present the least coarse granularity. Subsequent aggregation then zeroes in on the next flow key, now the lowest-ranked of the survivors. The records emanating from this second aggregation tier are stationed in buffer $B_2$, presenting a marginally coarser granularity view.

This iterative process persists through subsequent flow keys, ushering in an incremental level of flow granularity with every pass. The culmination is the aggregation of the highest-ranking flow key, where the consequent flow records—residing in buffer $B_4$—capture the coarsest granularity snapshot.

In summation, the Gradual Flow Aggregation strategy adeptly manages data granularity, striking a harmonious balance between storage efficiency and insightful flow data richness.

### B. Enhanced Gradual Flow Aggregation: A Refined Approach

Our proposed methodology optimizes the original Gradual Flow Aggregation scheme, introducing structure, adaptability, and improved data management. Central to this refinement is a multi-buffer structure reminiscent of the original method. However, what sets our method apart is its disciplined, sequential flow aggregation. Contrary to the original method, which allowed flow records the potential freedom to transition directly from the primary buffer to the last, our approach ensures flow records aggregate over only one flow key element when transitioning between consecutive buffers. This mandates a smoother, less abrupt data reduction, leading to enhanced information preservation.

Our methodology also integrates an improved threshold mechanism. Each buffer possesses its specific threshold, which can be consistent or varied across buffers. This threshold serves as a deciding factor, determining the flows to retain or those designated for aggregation. Moreover, by adjusting this mechanism to be responsive to resource utilization, we achieve a delicate balance between preserving pertinent flow details and efficiently managing resources.
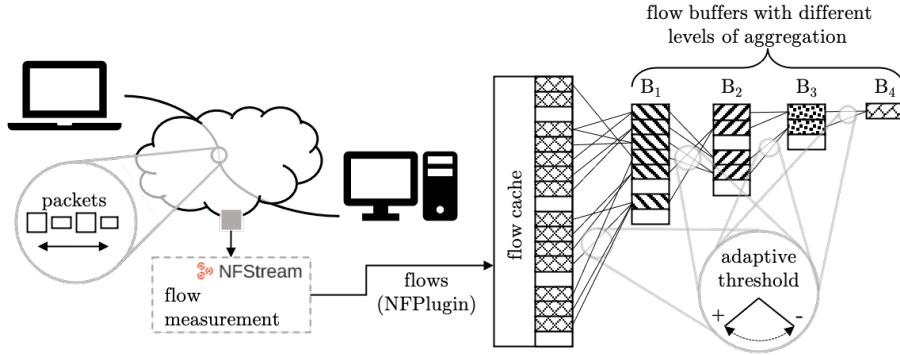
Fig. 2. Workflow from flow measurement achieved by NFStream to adaptive gradual flow aggregation implemented as an NFPlugin.

A notable innovation is the provision to redirect overflowing flows to subsequent buffers when a buffer reaches its capacity. This ensures that the aggregation process remains structured, even under significant network traffic demands. Furthermore, the last buffer, B4, is designed to capture all possible network flows, thus offering a holistic aggregation framework.

In summary, our enhanced methodology, characterized by its tiered flow key element reduction and dynamic threshold mechanism, not only minimizes the risk of sudden data reductions but also provides a more resource-efficient, organized, and information-rich flow aggregation process compared to its predecessor.

### C. Proof-of-Concept Implementation

Our methodology for enhanced gradual flow aggregation was built upon the NFStream framework [26], a Python-based tool engineered for rapid, flexible, and expressive data handling.

NFStream's robustness stems from its hybrid design, blending Python's accessibility with the speed of C. At the heart of this design is the NFlow structure, defined in C. This structure represents a network flow, encapsulating core attributes such as the five-tuple flow key, and is further enriched with additional flow statistics and metadata. These metadata are organized into multiple categories, including core features, L7 visibility, post-mortem statistics, and SPLT details. Users can dynamically toggle these feature sets on or off during flow measurement, offering a bespoke network analysis environment. Each of these features is stored in a specific variable type, tailored to accommodate its potential maximum value size.

NFStream not only simplifies the transition from raw network measurements to refined data science analytics due to its core functions in flow measurement and feature computation but also shines with its extendable architecture. This extensibility, particularly evident in the integration of custom network functionalities via the NFPlugin component [26], allowed us to embed our adaptive gradual flow aggregation directly as an NFPlugin. This integration was facilitated by the use of four distinct buffers ($B_1$-$B_4$), which collaboratively operate alongside NFStream's primary flow cache, as illustrated in Figure 2. These buffers are designed to retain flow records with varying granularities.

### D. Resource Utilization-driven Adaptability

Our approach to flow aggregation is adaptive, with the degree of aggregation varying in accordance with changes in resource utilization by the flow meter. The number of processed flows has been found to be directly proportional to resource utilization [27]. Essentially, an increase in network traffic, leading to a growing number of flows, results in a concurrent rise in CPU and memory load. Conversely, a reduction in traffic, and thus the number of flows, leads to a decrease in resource utilization. Assuming that the temporal aspect of the flow meter's resource utilization aligns with its flow cache load, we designed our aggregation process to reflect the changes of its CPU and memory usage.

With the above in mind, our methodology involves ongoing monitoring of CPU and memory utilization. This yields a set of $n$ observations denoted as $\mathcal{O} = o_t, o_{t+1}, \ldots, o_{t+n}$ ($t = 0, 1, \ldots, k$). Each observation is the average CPU and memory usage at time $t$, or $o_t = (cpu_t + memory_t)/2$. CPU time and memory usage are combined in our methodology to provide a more comprehensive measure. Analyzing these two metrics separately might overlook scenarios where one resource is heavily utilized while the other is not, resulting in a misleading representation of the overall system load. By combining these metrics, a more holistic and accurate assessment of the current state of system resource utilization is obtained. Furthermore, each observation, $o_t$, is expressed as a percentage, representing the proportion of the total available CPU and memory resources that are currently being used. Therefore, the value of $o_t$ in our approach can range from 0 to 100%.

The average resource utilization, represented as $\sigma$, is then calculated using the following formula:

$$\sigma = \begin{cases} \dfrac{\sum\limits_{i=k-w}^{k} o_{t+i}}{w}, & \text{if} \quad k \geq w. \\ \dfrac{\sum\limits_{i=0}^{k} o_{t+i}}{k}, & \text{otherwise.} \end{cases} \quad (1)$$

The sliding window of size $w$ is employed to ensure that our assessment of average resource utilization is up-to-date and sensitive to recent changes, by considering only the last

$w$ observations. However, at the beginning of our monitoring process, it may take some time to accumulate $w$ observations. The second case in the equation accounts for this scenario, where the average resource utilization is calculated using all available observations, instead of just the last $w$. Once we have collected at least $w$ observations, the first case will be used to calculate the average resource utilization, effectively implementing the sliding window approach.

Achieving adaptability involves comparing average resource utilization ($\sigma$) between consecutive times, $t-1$ and $t$. The aggregation threshold is adjusted based on the difference between $\sigma_{t-1}$ and $\sigma_t$. This can be formally expressed using the following formula:

$$T_{t+1} = \begin{cases} T_t + (|\sigma_{t-1} - \sigma_t|)\% \text{ of } T_{t-1}, & \text{if } \sigma_{t-1} < \sigma_t. \\ T_t - (|\sigma_{t-1} - \sigma_t|)\% \text{ of } T_{t-1}, & \text{if } \sigma_{t-1} > \sigma_t. \\ \varnothing, & \text{otherwise.} \end{cases}$$
(2)

where

$T_{t+1}$ = signifies the aggregation threshold for time $t+1$;

$T_t$ = represents the currently employed aggregation threshold at time $t$;

$T_{t-1}$ = refers to the aggregation threshold used at time $t-1$;

$\sigma_{t-1}$ = designates the previous average resource utilization as measured by CPU and memory load at time $t-1$;

$\sigma_t$ = indicates the current average resource utilization as measured by CPU and memory load at time $t$.

Equation (2) ensures that the aggregation threshold is regularly fine-tuned in accordance with variations in traffic load as manifested in CPU and memory utilization. As such, resource utilization near its maximum will trigger more aggressive flow aggregation, while lower resource utilization will result in a lesser degree of flow aggregation.

The degree of aggregation is based on a threshold $T$. Flows that transfer bytes ($flow_{bytes}$) equal to or exceeding $T$ are preserved in their original state. Conversely, flows that transfer fewer bytes than what $T$ prescribes are gradually aggregated within the buffers.

Aggregation commences once a buffer reaches its maximum capacity. Prior to adding a flow entry into the flow cache, the current buffer's capacity is assessed. Should it be at full capacity, the aggregation procedure is invoked. Flows adhering to the adaptively adjusted threshold remain intact within their existing buffer. Conversely, flows falling short of the threshold undergo aggregation. This entails discarding the flow key of the underlying buffer before relocating the flows to said buffer. This process is consistent across all buffers: each time a buffer reaches its limit, aggregation ensues, relegating flows to subsequent buffers characterized by diminished information retention capacities.

### E. Assumptions

Our experimental evaluation is anchored on several foundational assumptions.

1) *Hierarchy of Flow Features*: Building upon the guidance of [8] and [17], we delineated a hierarchy for flow features, sequenced from highest to lowest precedence as: *protocol*, *source port*, *destination port*, *source IP*, and *destination IP*. This structure aims to replicate a genuine scenario, ensuring that flows in the concluding buffer ($B_4$) can be differentiated based purely on their protocol identifiers.

2) *Aggregation and Flow Movement*: We postulate that aggregation happens solely within each discrete buffer, with the flow cache being an exception. Flow records transition from the primary flow cache towards the buffer $B_4$, as illustrated in Figure 2. This directional progression of records guarantees that transporting a record from the flow cache to $B_1$ is more cost-efficient than shifting records between buffers $B_3$ and $B_4$. This movement approach is optimized to curtail information loss, especially for flow records with significant byte transfer. To further impede information loss during synchronization, a flow record's relocation is restricted to one move per aggregation cycle.

3) *Adaptive Thresholding*: The adaptation mechanism hinges on the threshold $T$, dictated by the overall data volume in the flow, denoted as $flow_{bytes}$. Our prototype harnesses a universal threshold influencing all buffers, encompassing the flow cache. Future studies should probe into the advantages of individual thresholds for each buffer.

4) *Sliding Window Size*: For our prototype evaluation, we fixed the sliding window size at 10 elements. This choice is underpinned by thorough examination of diverse settings and the resultant outcomes. Yet, more investigation is warranted to gain deeper insights into the ramifications of varied window sizes.

5) *Memory Consumption and Flow Record Size*: Our results are simplified with the assumption that each flow record consumes 21 bytes, given by the size of the flow key. The memory consumption of flow keys in NFStream are:

- Source & Destination IPv4 Address: 8 bytes each
- Source & Destination IPv6 Address: 16 bytes each
- Source & Destination Port: 2 bytes each
- Protocol Identifier: 1 byte

This culminates in a memory consumption of 21 bytes for IPv4 flow records and 37 bytes for IPv6 flow records. For instance, storing 1 million IPv4 flow records would require approximately 20.96 MB. Nonetheless, for actual memory allocations, the results should be multiplied by the size of the complete feature set, computed as the sum of all feature sizes. Factoring in the 86 flow features NFStream can measure amplifies the memory requirements, emphasizing the need for effective memory management in IP flow measurement.

6) *TCP Flow Record Aggregation*: During the aggregation phase for TCP flow records, TCP flags are omitted. This is because these flags, being intrinsically flow-specific, lose their informational significance upon aggregation, rendering them redundant in the aggregated context.

## F. Information Bias in Adaptive Flow Aggregation

Through the process of aggregation, the information value of flow records diminishes, and this effect intensifies with increasingly aggressive flow aggregation, a consequence of surging resource utilization. Nevertheless, the impact on traffic management-related activities is direct and substantial. Information bias, in this context, refers to the distortion of flow size data resulting from the aggregation process. As the flow records are aggregated and the granularity of the information reduces, the interpretation of these records can become less precise, or "biased".

Consider the varied sensitivities of different applications to this information bias. For example, application type classification and DoS attack detection are highly sensitive to information bias. In application type classification, a detailed breakdown of the flow is essential to accurately identifying the type of application from its network patterns. Similarly, DoS attack detection relies heavily on identifying anomalies or spikes in individual flows, which can be obfuscated by aggregated data.

On the other hand, certain applications might be less affected by this bias. Accounting applications, for instance, often focus on the overall data transfer volume, which can be accurately measured even with aggregated flow data. Load balancing, too, is typically concerned with overall network utilization across multiple routes or servers, rather than the details of individual flows, and thus may remain largely unaffected by the aggregation process.

Therefore, understanding the effects of adaptive gradual flow aggregation on the value of flow size information, particularly in relation to the saturation of flow cache size, is critical. It is a balancing act between maintaining network performance and retaining the precision necessary for certain network tasks. Quantifying how flow size-based aggregation across the buffers $B_1 - B_4$ contributes to information bias can provide key insights that improve network traffic management. It is worth noting that while 'information bias' may not be a universally recognized term in this context, it serves effectively to describe the phenomenon discussed here.

## G. Dataset Preparation and Ground Truth Establishment

Our study utilized the UNIV1 dataset [28], [29], a publicly accessible traffic trace collected from a university campus data center. The dataset encompasses a wide array of services, including system backups, distributed file system hosting, e-mail servers, web services, and multicast video streams.

Organization of packets into forward, backward, and bidirectional flows was carried out using NFStream, with the passive and active expiration of flows set to NFStream's default of 120 and 1800 seconds, respectively. The resultant dataset comprised a total of 468,905 IP flows distributed across 14 distinct protocols. Specifically, UDP, TCP, and other protocols (such as ICMP, IGMP, EGP, IGP, ESP, and AH) accounted for 270,028, 196,305, and 2,572 flows respectively.

The CDF of flow sizes is displayed in Figure 3. Detailed analysis revealed that around $85\%$ of flows were smaller than 10 kB, with the majority of the data (20 kB and larger)
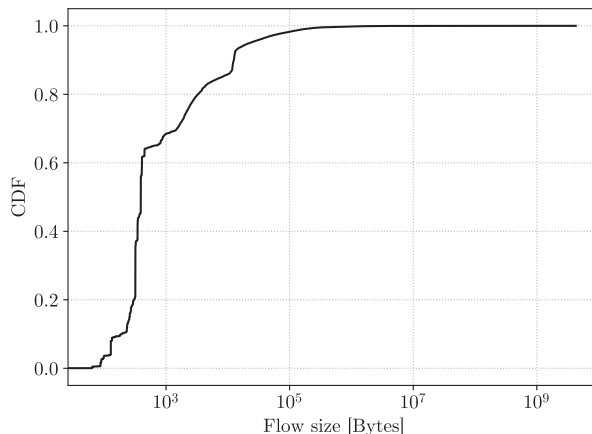


Fig. 3. Flow Size Distribution.

concentrated in the top $5\%$ of flows. Interestingly, less than 10 packets generated roughly $70\%$ of the flows, while a mere $2\%$ of all flows originated from more than 100 packets. We also discovered that about $14\%$ of all flows in the UNIV1 dataset transferred data larger than 10 kB in size. This observation is consistent with findings in [29], which reported that $10\%$ of the flows transport the majority of the traffic. The slight difference between our measurements and these findings can be attributed to variations in flow metering methodologies and the flow expiration timeout used.

Given our observations and the significant influence of flow size on our adaptive aggregation implementation, we selected the transmitted number of bytes in the flow ($flow_{bytes}$) as the primary metric for our evaluation. Consequently, we established the flow size threshold at 13.4 kB, which has been recognized as the optimal threshold for the UNIV1 dataset concerning flow size [30]. This threshold served as the benchmark for our evaluation.

## H. Quantifying Information Bias with Size Estimation Error

To study the impact of adaptive gradual flow aggregation on information bias, we introduce Size Estimation Error (SEE). The SEE is a measure of the error in the estimated flow size for flows within the buffers ($B_1$-$B_4$) at the conclusion of each measurement interval, relative to the memory size $m_s$.

*1) Definition of SEE:* SEE serves as an indicator of the average number of flows merged within the buffers and the relative change in their sizes compared to their original dimensions. For any given buffer $i$ ($B_i$), we calculate SEE as:

$$S\hat{E}E_i = |B_i| \times (\frac{1}{4} \times i). \tag{3}$$

In this equation, $S\hat{E}E_i$ is a preliminary estimation of the Size Estimation Error for buffer $i$. We employ the hat (^) notation to signify a preliminary value or estimate within the statistical framework.

*2) Normalization of SEE and Definition of SEE Vector:*
To acquire a normalized SEE ($SEE_i$) that provides a relative
estimation error for each buffer, we normalize each $\hat{SEE}_i$ as
follows:

$$SEE_i = \frac{\hat{SEE}_i}{\sum_{j=1}^{4} \hat{SEE}_j}. \tag{4}$$

Each $SEE_i$ now represents the normalized size estimation
error for buffer $i$, providing a proportionate share of the total
estimation error.

We then compile these normalized size estimation errors
into a single vector, referred to as the SEE vector:

$$SEE = (SEE_1, SEE_2, SEE_3, SEE_4). \tag{5}$$

We operate under the premise that the sum of all $SEE_i$
values is 1, and each $SEE_i$ is equal to or greater than 0, i.e.,
$\sum_{i=1}^{4} SEE_i = 1, \quad SEE_i \geq 0$.

*3) Adaptive Gradual Flow Aggregation and SEE Implemen-
tation:* Our implementation of adaptive gradual flow aggrega-
tion encompasses two potential scenarios for transferring flows
from a lower-ranked buffer $B_n$ to a higher-ranked buffer $B_{n+1}$.
In striving for a realistic environment, we set the memory limit
of each buffer to 10,000 bytes, aligning with the parameter
settings proposed by [14].

The first scenario comes into play when the memory reaches
saturation, necessitating memory clearance. Consequently, the
adaptive gradual flow aggregation relocates flows into lower-
ranked buffers, with a particular focus on moving the smallest-
sized flows between buffers.

In the second scenario, the system iteratively examines the
flows in buffers $B_1 - B_4$ based on the size of the $flow_{bytes}$
parameter (*cf.* Section IV-D). The flows are then tagged and
moved in accordance with the actual threshold $T$. As the
system nears the memory limit, it begins to reposition the
flow records relative to the current system utilization, adhering
to the adaptive gradual flow aggregation scheme discussed in
Section IV-A.

*4) Evaluation of Information Bias:* After distributing all
the flows across buffers $B_1 - B_4$, we contrast the results
to quantify the information bias induced by adaptive gradual
flow aggregation on the flow size feature. Importantly, our
method retains all necessary metadata for reverting aggregated
flows over the buffers back to their original form, enabling an
accurate efficacy evaluation.

## V. RESULTS

This section delineates the quantitative ramifications of
adaptive gradual flow aggregation. Initially, it assesses the
parameters of a system operating under constraints. Subse-
quently, it explicates the influence of adaptive flow aggregation
in a multi-buffer system arrangement.

### A. Bounds of a Constrained System Operation

Gauging the necessities and constraints for full-fledged,
dependable flow metering presents a formidable challenge. In
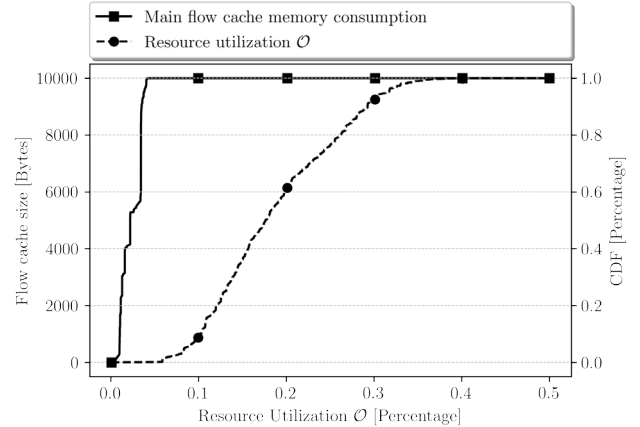spite of only a fraction of the traffic being measured, the



Fig. 4. Correlation between flow occupancy and resource utilization $\mathcal{O}$.

network invariably houses more devices than the metering
system. This culminates in a vast discrepancy between the
resources available for traffic generation and those for traffic
measurement, leading to a pronounced resource asymmetry. To
quantify this asymmetry, our initial step involved examining
the bounds for constrained operation of our prototype imple-
mentation devoid of active adaptive gradual flow aggregation.

Figure 4 visualizes the correlation between per-buffer flow
occupancy size and resource utilization, denoted as $\mathcal{O}$. The
left $y$-axis signifies buffer memory usage in bytes, the right
$y$-axis represents the distribution of flow sizes via CDF, and
the $x$-axis denotes resource utilization $\mathcal{O}$, which comprises
memory and CPU loads. The system's behaviour is depicted
in Figure 4 through two curves: (*i*) the distribution of flow
occupancies relative to $\mathcal{O}$, marked using rectangular markers,
and (*ii*) the CDF relative to $\mathcal{O}$, indicated using dot markers.

Figure 4 reveals that a surge in the number of flow records
in the flow cache directly leads to a significant increase
in resource utilization $\mathcal{O}$, which is composed of memory
and CPU loads. Consequently, the maximum memory cap of
10,000 bytes is attained at a resource utilization as low as
4.81%. As the observation of additional flows sustains the
memory load at 100%, the combined resource utilization $\mathcal{O}$
promptly escalates beyond 40%. This translates to both a
subpar operation of the system and the dropping or discarding
of packets and flows. We also noted that there is a 9.88%
chance of observing 10% $\mathcal{O}$ at a memory load of 1,420 bytes.

This implies that memory saturation, regardless of flow
creation and maintenance, substantially elevates the $\mathcal{O}$ param-
eter of a given device. Therefore, flow meters ought to be
calibrated with respect to memory limits to prevent reaching
the saturation level. In situations where such conditions are
unattainable, adaptive gradual flow aggregation can aid in
guaranteeing a reliable system operation.

### B. Implications of Adaptive Gradual Flow Aggregation

Adaptive gradual flow aggregation is conducted relative to
the actual memory load, capped at 10,000 bytes, and CPU
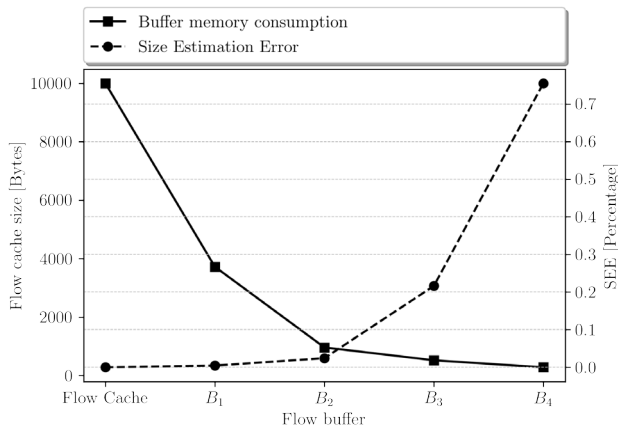utilization. This leads to the distribution of flows among indi-

Fig. 5. Correlation between the size of per-buffer flow occupancy and SEE.

vidual buffers. With this in mind, we explored SEE resulting from this distribution. Accordingly, packets in the UNIV1 were arranged using adaptive gradual flow aggregation, and the buffer status was examined post-operation. Note that, for the evaluation of this approach's operational conditions, flow expiration causing flows to be flushed from the cache(s) was disregarded.

Figure 5 illustrates the correlation between the size of per-buffer flow occupancy and SEE. The left $y$-axis signifies buffer memory usage in bytes, the right $y$-axis represents SEE, and the $x$-axis denotes the flow buffers. Figure 5 portrays the system behavior via two curves: ($i$) the distribution of flow occupancies between the buffers, marked using rectangular markers, and ($ii$) the SEE for each buffer, indicated using dot markers.

Figure 5 reveals that the buffers' occupancy by flows decreased logarithmically due to gradual flow reduction. The most substantial reduction occurred in buffer $B_4$, which aggregates flows over the flow key with the highest rank—the protocol identifier (cf. Section IV-A). Quantitatively, the SEE amounted to 75.49%, signifying a considerable information loss of all the UNIV1 flows. Nevertheless, this occurred at the cost of a significant reduction factor, implying all flows moved into this buffer were aggregated to yield 14 entries, each 21 bytes in size (UNIV1 contains 14 unique protocols, as discussed in Section IV-G).

Shifting towards a configuration devoid of adaptive gradual aggregation, the ratio discernibly shifts in favor of preserving information value, albeit at the expense of deteriorated memory utilization. Empirically, buffer $B_3$, which aggregates UNIV1 flows over the flow keys with the second-highest rank, yielded an SEE of 21.63%, consuming 540 bytes of buffer memory. Buffer $B_2$ presented a further improvement, with an SEE of 2.42% at a memory utilization of 960 bytes. Finally, aggregation over buffer $B_1$ resulted in an SEE of 0.46%, occupying 3,840 bytes of memory.

In light of the above, adaptive gradual flow aggregation over $B_2$ appears to offer the most optimal balance, providing an excellent compromise between preserving maximum informa-

tion value and minimal memory utilization, as demonstrated in Figure 5. Buffer $B_2$ aggregates flow records over the flow key with the second-lowest rank among all flow keys, thereby the information contained in the aggregated flow records has a less fine level of granularity compared to complete flow records, inevitably leading to information loss. Quantitatively, this approach aggregated 20,434 flows in $B_2$ from a total of 468,905 flows in UNIV1, translating to a loss of 2.42%. This could be considered a reasonable trade-off for enhanced resource utilization. Furthermore, the extent of loss can be mitigated by bolstering the physical hardware, specifically the allocated memory. Nevertheless, adaptive gradual flow aggregation significantly improved the system operation under limited constraints. It can aid in achieving optimal resource utilization while retaining the maximum information value embedded in the flow records.

## VI. Discussion

The adaptive gradual flow aggregation mechanism presented in this research study aims to optimize the balance between memory utilization and information value preservation in the context of network flow metering. The results shed light on the implications of adopting this mechanism, particularly in scenarios where the system is resource-constrained.

In a context where there exists a significant asymmetry between traffic generation and traffic measurement resources, our adaptive flow aggregation technique comes into play. It addresses the problem by efficiently distributing flows across multiple buffers, a fact underpinned by our results that highlight the performance of the mechanism under constrained operation.

Equally important is the role of SEE, a measure introduced to quantify the error in flow size estimation. Through the use of SEE, we can understand and assess the trade-off between memory utilization and information preservation. Indeed, our findings show that, as the level of flow aggregation increases, SEE increases too, which translates to greater information loss. However, it is noteworthy that the amount of information loss varies significantly across buffers, with some buffers offering an optimal trade-off.

Our results revealed that saturation, irrespective of the occurrence of flow creation and maintenance, causes the resource utilization parameter $\mathcal{O}$ of a given device to rise significantly. This underlines the importance of appropriately configuring flow meters concerning memory bounds to prevent reaching the saturation level. In scenarios where this is not feasible, adaptive gradual flow aggregation can ensure reliable system operation.

It was also observed that the application of adaptive flow aggregation resulted in a reduction in the total number of flow records in the cache, with a corresponding decrease in the SEE. The interplay between memory utilization and information preservation became apparent, illustrating the trade-offs involved. The buffer occupancy, and consequently the SEE, decreased logarithmically due to gradual flow reduction. While this led to a loss in information, it also enhanced system performance by significantly reducing memory usage.

An intriguing finding was the optimal balance achieved with adaptive gradual flow aggregation over buffer $B_2$. This offered an excellent trade-off between maximum information value and minimal memory utilization. Buffer $B_2$ aggregates flow records over the flow key with the second-lowest ranking among all flow keys, implying the information contained in the aggregated flow records has a coarser level of granularity compared to complete flow records. This led to an inevitable, yet acceptable, information loss. Nevertheless, with enhanced physical hardware or increased allocated memory, the degree of this loss can be minimized, hinting at potential avenues for improvement.

Overall, the results demonstrate the effectiveness and potential of our proposed adaptive gradual flow aggregation mechanism. Despite the inherent information loss, the mechanism significantly improves system operation under limited constraints and resource asymmetry. It achieves optimal resource utilization while striving to retain the maximum possible information value in the flow records. As such, it presents a promising solution for environments constrained by memory and processing resources.

## VII. IMPLICATIONS

The absence of a systematic data reduction strategy, while ensuring full preservation of information value, incurs higher operational costs and risks a potential crash of the measurement instrument as resource capacity reaches its limits. On the other hand, the application of gradual flow aggregation reduces the size of flow measurement data, thus lowering operational costs, but inevitably distorts flow-level information. Our study reveals that the operation of such a system can be further optimized for better resource utilization through adaptability.

From a traditional viewpoint, achieving optimal operation without adaptive flow aggregation seems impracticable as the dual objectives of preserving information value and reducing data volume are essentially mutually exclusive. Moreover, optimizing for either of these objectives leads to increased operational costs. In contrast, adaptive flow aggregation enables system operation optimization by balancing between information preservation and data volume reduction. The goal of optimization, therefore, is to determine the aggregation threshold for a given buffer $n$, given factors such as flow record creation rate, buffer capacity, memory utilization, and CPU utilization. The aim is to minimize the volume of flow records while maximizing the preservation of information in the flow records.

Our results show that optimizing system operation through preliminary observation-based configurations can already yield a more compact data volume while preserving considerable information value, all without incurring additional operational costs. This operation can be enhanced by solving the optimization problem to determine the optimal threshold per buffer $n$, including the flow cache. However, further research is necessary to establish the efficacy of algorithms designed to solve this optimization problem. These investigations could provide additional insight into the interplay between information preservation, data volume reduction, and resource utilization in flow metering systems.

## VIII. BROADER APPLICATIONS AND FUTURE DIRECTIONS

Given the versatility and profound implications of network monitoring, our study extends beyond merely reducing data volume to encompass a broad spectrum of applications:

*1) Quality of Service:* Effective flow aggregation is instrumental in accurately monitoring service levels. By ensuring network services consistently meet their designated quality parameters, we can enhance SLA adherence and elevate the overall network experience for users.

*2) IoT and Edge Computing:* The inherent constraints of IoT devices and the trend towards edge computing highlight the indispensability of adaptive flow aggregation. In these settings, striking a balance between data richness and resource efficiency is of paramount importance.

*3) Software-Defined Networking:* In the realm of SDNs, it is essential to fine-tune flow table rules and guarantee efficient load distribution. Our adaptive aggregation approach offers a promising avenue for honing rule sets and optimizing traffic management.

*4) Attack Detection:* Leveraging adaptive flow aggregation allows for the efficient clustering of disparate anomalous patterns. This capability can significantly bolster intrusion detection systems and fortify network security frameworks.

Peering further into areas like attack detection and QoS, we discern the potential for synergizing our adaptive flow aggregation techniques with advancements in AI and machine learning. Such a fusion could usher in predictive flow management, seamlessly navigating the dichotomy between preserving information and reducing data volume. Moreover, we envision subsequent research endeavors delving into various contexts and proposing algorithmic solutions to the optimization challenge we have underscored. Crafting these solutions would enhance the potency of adaptive flow aggregation, optimizing flow metering systems, especially in resource-limited environments.

## IX. CONCLUSION

In this study, we examined adaptive gradual flow aggregation as a methodology for coping with resource asymmetry in flow metering systems. Our results illuminated the practical implications of adaptive flow aggregation and highlighted the inherent trade-off between memory utilization, CPU load, and information preservation. By exploiting this trade-off, we found that adaptive gradual flow aggregation could help achieve a desirable balance between resource utilization and the maintenance of maximum information value within the flow records. In particular, we identified that buffer $B_2$ offered the most optimal balance, yielding a compromise loss of $2.42\%$, which we consider acceptable in the context of improved resource utilization.

Overall, our study underscores the potential of adaptive gradual flow aggregation to improve resource utilization in flow metering systems while preserving vital flow-level information. Adhering to open science principles, we have made the scripts used in our experiments publicly accessible[1]. We

---

[1] https://github.com/FlowFrontiers/AGFA

believe this will facilitate a thorough comprehension of our methodologies and encourage additional investigation in this domain.

## REFERENCES

[1] B. Li et al., "A survey of network flow applications," *Journal of Network and Computer Applications*, vol. 36, no. 2, pp. 567–581, 2013. DOI: 10.1016/j.jnca.2012.12.020.

[2] R. Hofstede et al., "Flow monitoring explained: From packet capture to data analysis with netflow and ipfix," *IEEE Communications Surveys Tutorials*, vol. 16, no. 4, pp. 2037–2064, 2014. DOI: 10.1109/COMST.2014.2321898.

[3] S. Bauer et al., "On the evolution of internet flow characteristics," in *Proceedings of the Applied Networking Research Workshop, ser.* ANRW '21, Proceedings of the Applied Networking Research Workshop, 2021, pp. 29–35. DOI: 10.1145/3472305.3472321.

[4] P. Velan, "Improving network flow definition: Formalization and applicability," in *NOMS 2018 – 2018 IEEE/IFIP Network Operations and Management Symposium*, NOMS 2018 – 2018 IEEE/IFIP Network Operations and Management Symposium, 2018, pp. 1–5. DOI: 10.1109/NOMS.2018.8406203.

[5] A. Dainotti, A. Pescape, and K. C. Claffy, "Issues and future directions in traffic classification," *IEEE Network*, vol. 26, no. 1, pp. 35–40, 2012. DOI: 10.1109/MNET.2012.6135854.

[6] S. Lee, K. Levanti, and H. S. Kim, "Network monitoring: Present and future," *Computer Networks*, vol. 65, pp. 84–98, 2014. DOI: 10.1016/j.comnet.2014.03.007.

[7] S. Dong and Y. Xia, "Network traffic identification in packet sampling environment," *Digital Communications and Networks*, 2022. DOI: 10.1016/j.dcan.2022.02.003.

[8] H. Irino, M. Katayama, and S. Chaki, "Study of adaptive aggregation on ipfix," in *Proceedings of the 7th Asia-Pacific Symposium on Information and Telecommunication Technologies*, Proceedings of the 7th Asia-Pacific Symposium on Information and Telecommunication Technologies, 2008, pp. 86–91. DOI: 10.1109/APSITT.2008.4653545.

[9] K. Cho, R. Kaizaki, and A. Kato, "Aguri: An aggregation-based traffic profiler," in *Quality of Future Internet Services*, M. I. Smirnov et al., Eds., Quality of Future Internet Services, 2001, pp. 222–242. DOI: 10.1007/3-540-45412-8_16.

[10] L. Kong et al., "Edge-computing-driven internet of things: A survey," *ACM Comput. Surv.*, 2022, Just Accepted. DOI: 10.1145/3555308.

[11] K. Cho, R. Kaizaki, and A. Kato, "An aggregation technique for traffic monitoring," in *Proceedings 2002 Symposium on Applications and the Internet (SAINT) Workshops*, 2002, pp. 74–81. DOI: 10.1109/SAINTW.2002.994556.

[12] G. Cheng and J. Gong, "Adaptive aggregation flow measurement on high speed links," in *Proceedings of the 11th IEEE Singapore International Conference on Communication Systems (ICCS)*, 2008, pp. 559–563. DOI: 10.1109/ICCS.2008.4737246.

[13] K.-c. Lan and J. Heidemann, "A measurement study of correlations of internet flow characteristics," *Computer Networks*, vol. 50, no. 1, pp. 46–62, 2006. DOI: 10.1016/j.comnet.2005.02.008.

[14] V. Sivaraman et al., "Heavy-hitter detection entirely in the data plane," in *Proceedings of the Symposium on SDN Research*, ser. SOSR '17, Proceedings of the Symposium on SDN Research, 2017, pp. 164–176. DOI: 10.1145/3050220.3063772.

[15] Y. Hu, D.-M. Chiu, and J.-S. Lui, "Adaptive flow aggregation – a new solution for robust flow monitoring under security attacks," in *10th IEEE/IFIP Conference on Network Operations and Management Symposium*, ser. NOMS '06, 10th IEEE/IFIP Conference on Network Operations and Management Symposium, 2006, pp. 424–435. DOI: 10.1109/NOMS.2006.1687572.

[16] Y. Hu, D. M. Chiu, and J. C. S. Lui, "Entropy based adaptive flow aggregation," *IEEE/ACM Transactions on Networking*, vol. 17, no. 3, pp. 698–711, 2009. DOI: 10.1109/TNET.2008.2002560.

[17] A. Pekar et al., "Adaptive aggregation of flow records," *Computing and Informatics*, vol. 37, no. 1, pp. 142–164, 2018. DOI: 10.4149/cai\_2018\_1\_142.

[18] C. Bi et al., "On precision and scalability of elephant flow detection in data center with SDN," in *Proc. 32nd IEEE Global Communications Conf. Workshops*, ser. GLOBECOM'13, 2013, pp. 1227–1232. DOI: 10.1109/GLOCOMW.2013.6825161.

[19] S. Wang et al., "Fdalb: Flow distribution aware load balancing for datacenter networks," in *2016 IEEE/ACM 24th International Symposium on Quality of Service (IWQoS)*, 2016, pp. 1–2. DOI: 10.1109/IWQoS.2016.7590409.

[20] S. Wang et al., "Flow distribution-aware load balancing for the datacenter," *Computer Communications*, vol. 106, pp. 136–146, 2017. DOI: 10.1016/j.comcom.2017.03.005.

[21] Z. Liu et al., "An adaptive approach for elephant flow detection with the rapidly changing traffic in data center network," *Int. J. of Network Management*, vol. 27, no. 6, e1987, 2017, e1987 nem.1987. DOI: 10.1002/nem.1987.

[22] N. Saha, S. Misra, and S. Bera, "Qos-aware adaptive flow-rule aggregation in software-defined iot," in *2018 IEEE Global Communications Conference (GLOBECOM)*, 2018, pp. 206–212. DOI: 10.1109/GLOCOM.2018.8647471.

[23] Q. T. Minh et al., "Flow aggregation for sdn-based delay-insensitive traffic control in mobile core networks," *IET Communications*, vol. 13, no. 8, pp. 1051–1060, 2019. DOI: 10.1049/iet-com.2018.5194.

[24] T. V. Phan et al., "Destination-aware adaptive traffic flow rule aggregation in software-defined networks," in *2019 International Conference on Networked Systems (NetSys)*, 2019, pp. 1–6. DOI: 10.1109/NetSys.2019.8854510.

[25] W.-K. Jia and X. Wang, "Flow aggregation for large-scale sdns with scattered address space allocation," *Journal of Network and Computer Applications*, vol. 169, p. 102 787, 2020. DOI: 10.1016/j.jnca.2020.102787.

[26] Z. Aouini and A. Pekar, "Nfstream: A flexible network data analysis framework," *Computer Networks*, vol. 204, p. 108 719, 2022. DOI: 10.1016/j.comnet.2021.108719.

[27] Y. Fu et al., "Jellyfish: Locality-sensitive subflow sketching," in *IEEE INFOCOM 2021 - IEEE Conference on Computer Communications*, IEEE, 2021, pp. 1–10. DOI: 10.1109/INFOCOM42981.2021.9488847.

[28] T. Benson, *Data set for IMC 2010 data center measurement*, University of Wisconsin-Madison, 2010. https://pages.cs.wisc.edu/~tbenson/IMC10_Data.html.

[29] T. Benson, A. Akella, and D. A. Maltz, "Network traffic characteristics of data centers in the wild," in *Proc. 10th Internet Measurement Conf.*, ser. IMC '10, Proc. 10th Internet Measurement Conf., 2010, pp. 267–280. DOI: 10.1145/1879141.1879175.

[30] A. Pekar et al., "Knowledge discovery: Can it shed new light on threshold definition for heavy-hitter detection?" *Journal of Network and Systems Management*, vol. 29, no. 3, p. 24, 2021. DOI: 10.1007/s10922-021-09593-w.

Balancing Information Preservation and Data Volume Reduction:
Adaptive Flow Aggregation in Flow Metering Systems

**Adrian Pekar** received the Ph.D. degree in computer science from the Technical University of Košice, Slovakia, in 2014. Currently, he is a Senior Researcher with the Department of Networked Systems and Services, Budapest University of Technology and Economics, Hungary. Prior to this, he held research, teaching, and engineering positions in Slovakia and New Zealand. His research interests include network and services management, software-defined networking, network function virtualization, and cloud computing.

**Laszlo A. Makara** earned his MSc degree from the Department of Networked Systems and Services at the Budapest University of Technology and Economics, Hungary, in 2023. Presently, he is pursuing his PhD at the same department, delving deeper into the world of computer engineering and research. His research is primarily focused on network and services management, software-defined networking, and network programmability.

**Winston K. G. Seah** received the Dr. Eng. degree from Kyoto University, Kyoto, Japan, in 1997. He is currently a Professor of network engineering with the School of Engineering and Computer Science, Victoria University of Wellington, New Zealand. Prior to this, he has worked for more than 16 years in mission-oriented industrial research, taking ideas from theory to prototypes, most recently, as a Senior Scientist with the Institute for Infocomm Research, Singapore. He has been actively involved in research in the areas of mobile ad hoc and sensor networks and co-developed one of the first QoS models for mobile ad hoc networks. His latest research interests include IoT, mobile edge computing, SDN, network anomaly detection, and 5G ultra reliable low latency and machine- type communications.

**Oscar M. Caicedo Rendon** (GS'11-M'15-SM'20) is a full professor at the University of Cauca, Colombia, where he is a member of the Telematics Engineering Group. He received his Ph.D. degree in computer science (2015) from the Federal University of Rio Grande do Sul, Brazil, and his M.Sc. in telematics engineering (2006) and his degree in electronics and telecommunications engineering (2001) from the University of Cauca. His research interests include network and service management, NFV, SDN, and Machine Learning for Networking.

# 5th Workshop on Management for Industry 4.0 – MFI5.0
## IEEE/IFIP Network Operations and Management Symposium 2024
## 6-10 May 2024 | Seoul, South Korea

**ORGANIZING COMMITTEE**

**Honorary Chair:**

Jerker Delsing
*LTU, SE*

**Organizing Chair:**

Markus Tauber
*RSA FG, AT*

**Technical Program Chair:**

Hans-Peter Bernhard *SAL & JKU Linz, AT*

**Technical Program Committee:**

Robert Harrison *Warwik University, UK*

Matthias Hemmje *FernUniversität, GER*

Nicholas Race *Lancaster University, UK*

Martin Wollschlaeger TU-Dresden, GER

Jürgen Jasperneite *Fraunhofer IOSB-INA & TH OWL, GER*

**IMPORTANT DATES**

Submission deadline:
January 19th, 2024

Author notification:
March 1st, 2024

Final submission:
March 15th, 2024

**QUESTIONS?**

markus.tauber@
researchstudios.at
http://mfi50.icb.at

**DESCRIPTION**

The transition from ISA-95 to RAMI4.0/IIRA based automation for production automation in Industry 4.0 is ongoing. This includes the integration of legacy OT with emerging IT technologies. Another aspect is automation/digitalization across value networks involving a multitude of stakeholders in complex relationships. Consequently, Management for Industry 5.0 covers three thematic themes: (1) connectivity, infrastructure, and security, (2) the autonomous evolution and challenges of System of System (SoS) in cyber-physical systems (CPS), and (3) the human in the loop.

Recent advancements in communication technology, especially wireless, are transforming the industrial landscape. This necessitates the integration of wireless and cellular tech, including 5G/6G, into both OT and IT communications, offering greater flexibility and challenging traditional industrial communication and security paradigms. This shift extends beyond industry 5.0, benefiting fields like agriculture and logistics. Key requirements such as privacy, dependability, and trustworthiness drive service- and data-driven automation in various production domains.

The future envisions large System of Systems (SoS) involving IoT, AI, Analytics, Big data, and legacy tech, distributed among multiple stakeholders. The success of these production systems hinges on incorporating human actors and addressing challenges like trust in autonomous systems, human-robot collaboration, competence development, and knowledge management. This encompasses explainable AI in production and workplace integrated learning in smart factories.

While architectures like RAMI4.0 and IIRA have been proposed, they are still in their early stages. Implementation platforms and frameworks are also in their infancy, particularly in managing complex automation and digitalization solutions across all levels. These architectures and technologies will be instrumental in autonomically controlling digitalized production infrastructures, requiring trustworthy and reliable data. Trust in industrial AI varies among stakeholders, impacting management and organizational aspects. Technology and organizational adaptation are critical, potentially leading to organizational reconfiguration.

**FOCUS ON:**

The workshop will focus on several core engineering and management issues, focus topics are:

- Migration Management
- Technology Acceptance
- Technology and Organizational Adaptation
- Operational Management
- Security Management
- Deployment Management
- Management of Networked Components in Industry 4.0/5.0 scenarios
- Technology and certification integration in CPS
- Automation evolution Management and Engineering

- Transition of I4.0 to I5.0
- Product Life Cycle Management
- Product Planning Management
- Manufacturing Change Management
- Manufacturing Process Management
- Manufacturing Operations Management
- Management of Digital Twins
- System of System Challenges in Cyber Physical System
- New Return on Investment and Sustainability Approaches

Additional topics may be considered given adequate proposal, therefore.

**SUBMISSION OF PAPERS:**

Authors are invited to submit original contributions written in English that have not been published or submitted for publication elsewhere. Technical papers must be formatted using the IEEE 2-column format and not exceed 6 pages for full paper submissions or not exceed 4 pages for short paper submissions. Papers should be submitted through NOMS submission system.

# Guidelines for our Authors

## Format of the manuscripts

Original manuscripts and final versions of papers should be submitted in IEEE format according to the formatting instructions available on

   *https://journals.ieeeauthorcenter.ieee.org/*
   *Then click: "IEEE Author Tools for Journals"*
   *- "Article Templates"*
   *- "Templates for Transactions".*

## Length of the manuscripts

The length of papers in the aforementioned format should be 6-8 journal pages.
Wherever appropriate, include 1-2 figures or tables per journal page.

## Paper structure

Papers should follow the standard structure, consisting of *Introduction* (the part of paper numbered by "1"), and *Conclusion* (the last numbered part) and several *Sections* in between.
The Introduction should introduce the topic, tell why the subject of the paper is important, summarize the state of the art with references to existing works and underline the main innovative results of the paper. The Introduction should conclude with outlining the structure of the paper.

## Accompanying parts

Papers should be accompanied by an *Abstract* and a few *Index Terms (Keywords)*. For the final version of accepted papers, please send the short cvs and *photos* of the authors as well.

## Authors

In the title of the paper, authors are listed in the order given in the submitted manuscript. Their full affiliations and e-mail addresses will be given in a footnote on the first page as shown in the template. No degrees or other titles of the authors are given. Memberships of IEEE, HTE and other professional societies will be indicated so please supply this information. When submitting the manuscript, one of the authors should be indicated as corresponding author providing his/her postal address, fax number and telephone number for eventual correspondence and communication with the Editorial Board.

## References

References should be listed at the end of the paper in the IEEE format, see below:
  a) Last name of author or authors and first name or initials, or name of organization
  b) Title of article in quotation marks
  c) Title of periodical in full and set in italics
  d) Volume, number, and, if available, part
  e) First and last pages of article
  f) Date of issue
  g) Document Object Identifier (DOI)

*[11] Boggs, S.A. and Fujimoto, N., "Techniques and instrumentation for measurement of transients in gas-insulated switchgear," IEEE Transactions on Electrical Installation, vol. ET-19, no. 2, pp.87–92, April 1984. DOI: 10.1109/TEI.1984.298778*

Format of a book reference:

*[26] Peck, R.B., Hanson, W.E., and Thornburn, T.H., Foundation Engineering, 2nd ed. New York: McGraw-Hill, 1972, pp.230–292.*

All references should be referred by the corresponding numbers in the text.

## Figures

Figures should be black-and-white, clear, and drawn by the authors. Do not use figures or pictures downloaded from the Internet. Figures and pictures should be submitted also as separate files. Captions are obligatory. Within the text, references should be made by figure numbers, e.g. "see Fig. 2."
When using figures from other printed materials, exact references and note on copyright should be included. Obtaining the copyright is the responsibility of authors.

## Contact address

Authors are requested to submit their papers electronically via the following portal address:

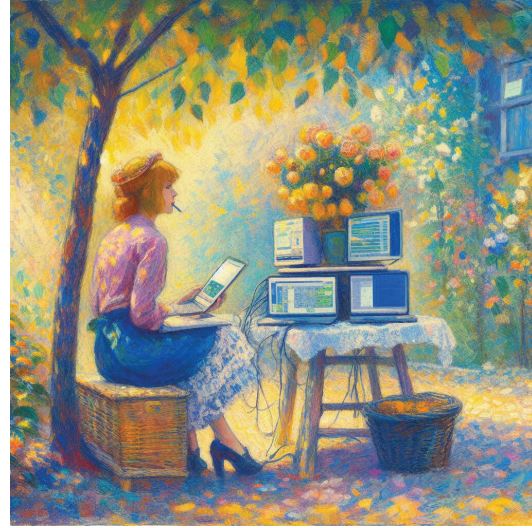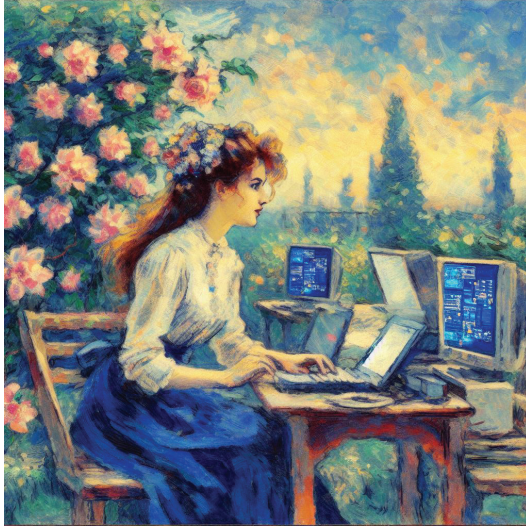https://www.ojs.hte.hu/infocommunications_journal/about/submissions

If you have any question about the journal or the submission process, please do not hesitate to contact us via e-mail:

Editor-in-Chief: Pál Varga – pvarga@tmit.bme.hu

Associate Editor-in-Chief:

József Bíró – biro@tmit.bme.hu

László Bacsárdi – bacsardi@hit.bme.hu

**2024 IEEE Network Operations and Management Symposium (NOMS 2024)**
**6-10 May 2024 Seoul, Korea**
*Towards intelligent, reliable, and sustainable network and service management*

**First IEEE Workshop on Generative AI for Network Management (GAIN)**

## • • • CALL FOR WORKSHOP PAPERS • • •

The GAIN workshop aims to systematically investigate and discuss the application of Generative AI in network management. It will bring together academic researchers from various disciplines (communication networks, data science, operational research) and practitioners from industry. Both scientific papers and industrial use case papers are welcomed. The considered topics in Generative AI for network management can initially be structured along the well-accepted FCAPS models in network management. This is a call for papers – example topics are as follows.

**Fault Management**
– Predictive Maintenance: Using generative AI for proactive network management
– Network Troubleshooting using Generative AI, incl. root cause analysis and resolution
– Monitoring using Generative AI: Using generative AI for efficient monitoring of network resources

**Configuration Management**
– Network Configuration Automation with Generative AI
– Automated Network Design and Deployment using Generative AI
– Generative AI for Traffic Management: Optimizing network traffic engineering through AI

**Accounting**
– Ethical Considerations: Addressing Privacy and Security Concerns in AI-Based Network Management
– Ensuring fairness between network users using generative AI for optimal resource allocation

**Performance**
– Network Optimization with AI: Leveraging generative algorithms for network efficiency
– Dynamic Resource Allocation: Leveraging generative AI for efficient network resource management
– Efficient Network Data Analysis using Generative AI:

**Security**
– AI-Based Security Protocols: Developing next-generation network security strategies
– Anomaly Detection and Response: Utilizing generative AI for enhanced network security

**Use Cases**
– Generative AI for management of IoT, Wireless/RAN, or Core, and Cloud-to-Edge networking

**General**
– Prompt Engineering for Network Management Using LLMs
– Robustness and Reliability of Generative AI for net. management (incl. benchmarks and datasets)
– Scalability Orchestration, Testing and Validation of Generative AI for Network Management

**Submission and Important Dates:**
Submission site: https://jems3.sbc.org.br/noms_gain2024
Paper Submission Deadline: **Jan. 19, 2024**
Notification of Acceptance: **Mar. 1, 2024**
Final Camera Ready: **Mar. 15, 2024**

**Workshop organisers:**
– Alberto Leon-Garcia, Univ. of Toronto, CA
  *(alberto.leongarcia@utoronto.ca)*
– Pal Varga, Budapest Univ. of Technology and Economics, HU
  *(pvarga@tmit.bme.hu)*
– Kurt Tutschku, Blekinge Inst. of Technology, SE
  *(ktt@bth.se)*

# SCIENTIFIC ASSOCIATION FOR INFOCOMMUNICATIONS



## Who we are

Founded in 1949, the Scientific Association for Info-communications (formerly known as Scientific Society for Telecommunications) is a voluntary and autonomous professional society of engineers and economists, researchers and businessmen, managers and educational, regulatory and other professionals working in the fields of telecommunications, broadcasting, electronics, information and media technologies in Hungary.

Besides its 1000 individual members, the Scientific Association for Infocommunications (in Hungarian: Hírközlési és Informatikai Tudományos Egyesület, HTE) has more than 60 corporate members as well. Among them there are large companies and small-and-medium enterprises with industrial, trade, service-providing, research and development activities, as well as educational institutions and research centers.

HTE is a Sister Society of the Institute of Electrical and Electronics Engineers, Inc. (IEEE) and the IEEE Communications Society.

## What we do

HTE has a broad range of activities that aim to promote the convergence of information and communication technologies and the deployment of synergic applications and services, to broaden the knowledge and skills of our members, to facilitate the exchange of ideas and experiences, as well as to integrate and harmonize the professional opinions and standpoints derived from various group interests and market dynamics.

To achieve these goals, we…

- contribute to the analysis of technical, economic, and social questions related to our field of competence, and forward the synthesized opinion of our experts to scientific, legislative, industrial and educational organizations and institutions;

- follow the national and international trends and results related to our field of competence, foster the professional and business relations between foreign and Hungarian companies and institutes;

- organize an extensive range of lectures, seminars, debates, conferences, exhibitions, company presentations, and club events in order to transfer and deploy scientific, technical and economic knowledge and skills;

- promote professional secondary and higher education and take active part in the development of professional education, teaching and training;

- establish and maintain relations with other domestic and foreign fellow associations, IEEE sister societies;

- award prizes for outstanding scientific, educational, managerial, commercial and/or societal activities and achievements in the fields of infocommunication.

## Contact information

President: **FERENC VÁGUJHELYI** • *elnok@hte.hu*
Secretary-General: **GÁBOR KOLLÁTH** • *kollath.gabor@hte.hu*
Operations Director: **PÉTER NAGY** • *nagy.peter@hte.hu*

Address: H-1051 Budapest, Bajcsy-Zsilinszky str. 12, HUNGARY, Room: 502
Phone: +36 1 353 1027
E-mail: *info@hte.hu*, Web: *www.hte.hu*