# Implementation and Experimental Evaluation Of a Virtualised 5G Network for URLLC Services

GIANNONE Francesco

# Implementation and Experimental Evaluation Of a Virtualised 5G Network for URLLC Services

## Author
Francesco Giannone

## Supervisor
Prof. Luca Valcarenghi

## Tutor
Prof. Piero Castoldi

# Contents

*Contents*

# List of Figures

*List of Figures*

# List of Tables

# 1 Introduction

5G Is Coming. Here is how immersive smartphones, autonomous cars and billions of intelligent sensors could transform your life.

The fifth generation of cellular technology, 5G, will offer much more than faster phones. Its first phase, scheduled to begin in 2019, will improve smartphone responsiveness. Soon after, 5G may also start to enhance many other aspects of your daily life, from automotive safety to your entertainment and your view of reality. For instance, your 5G smartphone will be able to stream ultra high-definition video at gi-



**Figure 1.1:** Orchestration in 5G

gabit speeds, with enough consistency to enable immersive virtual reality experiences.

In a subsequent phase, 5G could also provide a dependable link between autonomous vehicles, allowing cars to share their "intentions" and communicate directly with one another - and even with pedestrians - for safer, more efficient commutes.

5G, however, will not just be another "G" delivering faster speeds. Such standard brings technologies and facilitation that will allow to unite billions of devices, with the potential to transform not only your smartphone but also your home, vehicle, job, industry and community. Indeed, due to the inability in terms of power of our current generation of wireless internet devices to fully bring some services, such as virtual and augmented reality, into the daily life, 5G will be the only possible solution. That's because, for instance the extended reality - an umbrella term encompassing virtual reality, augmented reality and everything in between - pushes connectivity to its limits: it requires a

uniform high-speed connection, increased network capacity and minimal delays.



**Figure 1.2:** Ability to stream ultra high-definition video in 5G

5G will be able to support virtual teleconferencing. Say you are streaming a basketball game to your headset and receiving a 360-degree view of the court. However, you do not want to just watch the game; you also want to pause the action and brush shoulders with players. 5G will have the potential to provide enough capacity and low latency for you to immerse yourself in the game without delays or data congestion. Who wants to enter a virtual world, after all, only to be dragged out by slowed movements or scenery glitches? Extended reality enabled by 5G will also allow the users to virtually enter the same room as friends and family, with few delays or dropped connections. With the addition of haptic gloves - accessories that allow people to touch and feel in virtual reality - you might even be able to hold your child's hand while you are in another country.

Further, along in our 5G future, wearable technologies will join forces with the massive internet of things, IoT, allowing billions of intelligent, connected sensors to do more than tell you how many steps you've walked or how fast your heart is beating. Wearables of today provide limited information, but with the 5G massive IoT, they'll be able to connect to what's going on around us all of the time.



**Figure 1.3:** Future 5G Wearables

This is the 5G future: sensors could constantly interact with your wearable devices and with one another, creating the ability to monitor everything from your morning commutes to your factory goods to your children. A simple vibration from a wearable on your wrist, for instance, could tell you that only one bike remains at your usual bike-share location. As you rush over, the same

device might caution you to be aware of cars speeding through the nearest intersection.

5G massive IoT will have the ability to scale to billions of connected sensors. How would this play out? A chain of sensors could track and optimize goods as they move from their sources to your home. Take your clothes, for instance. At the beginning of the supply chain, sensors in the ground and in machinery could connect over the 5G network to farm cotton more efficiently. Self-driving trucks with sensors could then transport the raw cotton to a factory, where a wireless automated robot could transform the raw material into samples, based on size and pattern. Designers fitted with AR glasses could collaboratively turn the rough samples into templates for an entire production run. Finally, the shirts could be uploaded to an online shelf, where you could try them on using a VR device. This entire process, from farm to closet, would be brought to you by the 5G network.



**Figure 1.4:** 5G Connected Vehicles

5G will also enable ultra-reliable and low-latency communication. Autonomous vehicles and high-precision manufacturing, among other use cases, will need these qualities to thrive. 5G, at a later stage, could allow autonomous vehicles to establish a direct link between each other that bypasses network towers and diminishes urgent communications delays. This means a vehicle could instantly receive alerts to stop or change course if a collision seems imminent, even if the danger is beyond the range of its cameras and radar. In addition, 5G-enabled vehicles could connect with smart city infrastructure like traffic lights and street signs, improving traffic management. With 5G communications between autonomous vehicles, cars could share their intentions and know intended paths, allowing them to plan more safely and efficiently. Therefore, steering wheels could disappear, with cars redesigned as living spaces: basically couches or offices on wheels. With 5G-enhanced mobile internet, your car will be able to stream ultra high-definition video, and even VR, while you travel. Passengers could easily work as they ride, perhaps changing the entire nature of their commutes.

As the first 5G-enabled devices are released, 5G will be on the cusp of changing

the world in countless ways. And fully immersive virtual worlds, massive networks of connected sensors and safer roadways are just the beginning.

The imagination of our future is therefore a networked society with unbounded access to information and sharing of data which is accessible everywhere and ever time for everyone and everything. Present wireless based technologies, like 3GPP LTE technology, HSPA and Wi-Fi, will be incorporating new technology components that will be helping to meet the needs of the future. Nevertheless, there may be certain scenarios that cannot be adequately addressed along with the evolution of ongoing existing technologies. The instigation of a completely new wireless based technologies will complement the current technologies which are need for the long term realization of the networked society [1]

## 1.1 Novel contributions

This thesis provides contributions to the state of the art regarding the architecture and performance of the LTE and 5G mobile networks, from the conceptual and the experimental point of view. The thesis covers both research and implementation aspects, especially because almost everything was developed in an experimental setup to develop and analyze a standard compliant 5G network.

The main objective of the thesis is the development of a standard compliant 5G network and all the components belonging to it. In such a network, firstly a deep analysis of its performance have been performed. This is because the need to understand if some typical 5G services and applications, such as autonomous vehicles collision application and ultra high-definition video streaming applications could be deployed exploiting the deployed 5G network.

To achieve these goals, the starting point has been the investigation of one of the main open source 5G emulation platform, the OpenAirInterface emulation platform. After a first phase in which the platform has been analyzed and set in the TeCIP Institute laboratory, it has been extended with new features and capabilities and adopted to setup a standard compliant 5G network. A deeply analysis and experimental measurements of latency and jitter experienced in the 5G network have been performed.

The next objective was the introduction in the 5G network of the so-called Multi-access Edge Computing (MEC). Edge computing as an evolution of cloud computing brings application hosting from centralized data centres down to the network edge, closer to consumers and the data generated by applications. Edge computing

is acknowledged as one of the key pillars for meeting the demanding Key Performance Indicators (KPIs) of 5G, especially as far as low latency and bandwidth efficiency are concerned. This objective is mainly realized thanks to the collaboration with EURECOM in Sophia-Antipolis during the period abroad. Exploiting the MEC, two applications as part of our contribution to the 5G-TRANSFORMER project related to the road safety and automotive infotainment have been deployed in the autonomous vehicles scenario.

The 5G network architecture, will be heavily based on virtual network functions(VNFs). Thus, in a scenario where network functions are virtualized, both hardware and software failures assume the same importance, and their reliability shall be guaranteed. Similarly, reliability at service chain level is important to assure proper service availability features to application service platforms deployed by verticals. Therefore a study of the reliability in 5G networks has been also performed in the conclusion of my Ph.D. Different protection mechanisms have been proposed and experimentally analysed for both 5G core and access networks.

Last but not least, this work also pursued a really challenging goal: we were able to demonstrate a 5G network slice deployment in the ARNO testbed by using the 5G-TRANSFORMER architecture and offer a mobile/edge connectivity service with virtualized functions.

## 1.2 Thesis Structure

This Ph.D. thesis focuses on the design, implementation and experimental validation of a 5G network.

Chapter 2 gives a broad overview of the history of the mobile network starting from the first generation of mobile network (1G) until nowadays Long Term Evolution.

Chapter 3 gives a broad overview of the LTE architecture and key points and introduces the 5G Next Generation mobile network design and strong aspects. In this chapter, it is specified the difference between LTE and 5G network and the novelties introduced by the 5G.

Chapter 4 overviews the software platform and the hardware utilised in the TeCIP Laboratory a standard-compliant 5G network.

Chapter 5 presents an evaluation of how different virtualization technologies can afflict the network performances in terms of latency and jitter when different 5G network components are virtualised and different gNB functional splits are

considered. The measured latency and packet jitter (i.e., packet delay variation) "budget" are computed in all the possible combinations and compared with the scenario when all the network components belonging the 5G network are deployed in bare metal.

Chapter 6 illustrates the role and explains the benefits of Mobile-access Edge Computing introduction in a 5G network. Some possible solutions to deploy and integrate MEC in the 5G system are presented. Finally, it presents two deployed MEC-based applications in the automotive scenario.

Chapter 7 introduces and describes the reliability in the 5G networks, focusing on both the core network and access network. Firstly shows a proposed protection mechanism able to recover vEPC failures by means of a vEPC in "hot backup". Both working vEPC and backup vEPC are deployed in multiple Network Function Virtual Infrastructure Points of Presence (NFVI-PoPs) made available by the federated testbeds belonging to the SoftFIRE project. The goal is to evaluate the Service Recovery Time (SRT), that is the time required to regain user equipment (UE) connectivity, when the proposed resilient scheme is deployed in different NFVI-PoPs. Then, a Next Generation RAN testbed whose midhaul are realized using a highly reliable two layer Ethernet-over-DWDM transport network prototype is described. Finally, the application and experimental evaluation of two protection mechanisms to achieve reliable midhaul transport networks when using the functional split 7-1 option is presented.

Chapter 8 demonstrates a 5G network slice deployment in the ARNO test-bed by using the 5G-TRANSFORMER architecture and offer a mobile/edge connectivity service with virtualized functions.

Finally, Chapter 9 closes with a summary discussion, which includes an overview of the achieved results and contributions, as well as the list of scientific production of the Ph.D. period.

# 2  History of mobile network

## 2.1  Introduction

While the transmission of speech by radio has a long history, the first devices that were wireless, mobile, and also capable of connecting to the standard telephone network are much more recent. The first such devices were barely portable compared to today's compact hand-held devices, and their use was clumsy [2].

Along with the process of developing a more portable technology, and a better interconnections system, drastic changes have taken place in both the networking of wireless communication and the prevalence of its use, with smartphones becoming common globally and a growing proportion of Internet access now done via mobile broadband.

Before the devices existed that are now referred to as mobile phones or cell phones, there were some precursors. In 1908, a Professor Albert Jahnke and the Oakland Transcontinental Aerial Telephone and Power Company claimed to have developed a wireless telephone. They were accused of fraud and the charge was then dropped, but they do not seem to have proceeded with production. Beginning in 1918, the German railroad system tested wireless telephony on military trains between Berlin and Zossen. In 1924, public trials started with telephone connection on trains between Berlin and Hamburg. In 1925, the company Zugtelephonie AG was founded to supply train telephony equipment and, in 1926, telephone service in trains of the Deutsche Reichsbahn and the German mail service on the route between Hamburg and Berlin was approved and offered to first class travelers.

Fiction anticipated the development of real world mobile telephones. In 1906, the English caricaturist Lewis Baumer published a cartoon in Punch magazine entitled "Forecasts for 1907" in which he showed a man and a woman in London's Hyde Park each separately engaged in gambling and dating on wireless telephony equipment.

Then, in 1926, the artist Karl Arnold created a visionary cartoon about

the use of mobile phones in the street, in the picture "Wireless Telephony" (see Fig. 2.1), published in the German satirical magazine Simplicissimus. The Second World War made military use of radio telephony links. Hand-held radio transceivers have been available since the 1940s. Mobile telephones for automobiles became available from some telephone companies in the 1940s. Early devices were bulky, consumed large amounts of power, and the network supported only a few simultaneous conversations.

In the United States, engineers from Bell Labs began work on a system to allow mobile users to place and receive telephone calls from automobiles, lead-



**Figure 2.1:** Karl Arnold drawing of public use of mobile telephones

ing to the inauguration of mobile service on June 17th 1946 in St. Louis, Missouri. Shortly after, AT&T offered Mobile Telephone Service. A wide range of mostly incompatible mobile telephone services offered limited coverage area and only a few available channels in urban areas. The introduction of cellular technology, which allowed re-use of frequencies many times in small adjacent areas covered by relatively low powered transmitters, made widespread adoption of mobile telephones economically feasible.

In the USSR, Leonid Kupriyanovich, an engineer from Moscow, in 1957-1961 developed and presented a number of experimental pocket-sized communications radio. The weight of one model, presented in 1961, was only 70 g and could fit on a palm. However, in the USSR the decision at first to develop the system of the automobile "Altai" phone was made.

In 1965, the Bulgarian company "Radioelektronika" presented a mobile automatic phone combined with a base station at the Inforga-65 international exhibition in Moscow. Solutions of this phone were based on a system developed by Leonid Kupriyanovich. One base station, connected to one telephone wire line, could serve up to 15 customers.

The advances in mobile telephony can be traced in successive generations from the early "0G" services like Mobile Telephone Services (MTS) and its successor Improved Mobile Telephone Service (ITS), to first-generation (1G) analog cellular

**Figure 2.2:** Evolution of G-Technology

network, second-generation (2G) digital cellular networks, third-generation (3G) broadband data services, fourth-generation (4G) native-IP networks to the state-of-the-art five-generation (5G) Mobile Network. An evolution of the G-Technology is shown in Fig. 2.2.

## 2.2  1G - Analog cellular

The first type of mobile connectivity standard emerged in the 1970s, and lasted until the early '90s, at which point there were some 20 million subscribers worldwide. 1G was a pioneering technology that helped drive mass market usage of cellular technology, but it had several serious issues by modern standards. It was unencrypted and easily vulnerable to eavesdropping via a scanner; it was susceptible to cell phone "cloning" and it used a Frequency-division multiple access (FDMA) scheme and required significant amounts of wireless spectrum to support.



**Figure 2.3:** Motorola DynaTAC

The first analog cellular (1G) system widely deployed in North America was the

Advanced Mobile Phone System (AMPS). It was commercially introduced in the Americas in 13 October 1983, Israel in 1986, and Australia in 1987. The phone (see Fig. 2.3) had a talk time of just thirty-five minutes and took ten hours to charge.

## 2.3  2G - Digital cellular

In the 1990s, the "second generation" mobile phone systems emerged. Two systems competed for supremacy in the global market: the European developed GSM standard and the U.S. developed CDMA standard [3]. These differed from the previous generation by using digital instead of analog transmission, and also fast out-of-band phone-to-network signaling. The rise in mobile phone usage as a result of 2G was explosive and this era also saw the advent of prepaid mobile phones.

Second-generation 2G cellular networks were commercially launched on the GSM standard in Finland by Radiolinja (now part of Elisa Oyj) in 1991. Three primary benefits of 2G networks over their predecessors were that:

1. Phone conversations were digitally encrypted;

2. 2G systems were significantly more efficient on the spectrum enabling far greater wireless penetration levels;

3. 2G introduced data services for mobile, starting with SMS text messages, picture messages, and MMS (multimedia messages).

All text messages sent over 2G are digitally encrypted, allowing the transfer of data in such a way that only the intended receiver can receive and read it. With General Packet Radio Service (GPRS), 2G offered a theoretical maximum transfer speed of 50 kbit/s (40 kbit/s in practice) and with EDGE (Enhanced Data Rates for GSM Evolution), there was a theoretical maximum transfer speed of 1 Mbit/s (500 kbit/s in practice).



**Figure 2.4:** Personal Handy-phone System mobiles

The most common 2G technology was the time division multiple access (TDMA)-based GSM, originally from Europe but used in most of the world outside North America. Over 60 GSM operators were also using CDMA2000 in the 450 MHz frequency band (CDMA450) by 2010.

Coinciding with the introduction of 2G systems was a trend away from the larger "brick" phones toward tiny 100-200 grams hand-held devices as shown in Fig 2.4. This change was possible not only through technological improvements such as more advanced batteries and more energy-efficient electronics, but also because of the higher density of cell sites to accommodate increasing usage. The latter meant that the average distance transmission from phone to the base station shortened, leading to increased battery life while on the move.

## 2.4  3G - Mobile broadband

As the use of 2G phones became more widespread and people began to use mobile phones in their daily lives, it became clear that demand for data (such as access to browse the internet) was growing [4]. Further, experience from fixed broadband services showed there would also be an ever-increasing demand for greater data speeds. The 2G technology was nowhere near up to the job, so the industry began to work on the next generation of technology known as 3G. The main technological difference that distinguishes 3G technology from 2G technology is the use of packet switching rather than circuit switching for data transmission. In addition, the standardization process focused on requirements more than technology (2 Mbit/s maximum data rate indoors, 384 kbit/s outdoors, for example).

Inevitably this led to many competing standards with different contenders pushing their own technologies, and the vision of a single unified worldwide standard looked far from reality. The standard 2G CDMA networks became 3G compliant with the adoption of Revision A to EV-DO, which made several additions to the protocol while retaining backwards compatibility:

- Introduction of several new forward link data rates that increase the maximum burst rate from 2.45 Mbit/s to 3.1 Mbit/s;

- Protocols that would decrease connection establishment time;

- Ability for more than one mobile to share the same time slot;

- Introduction of QoS flags;

All these were put in place to allow for low latency, low bit rate communications such as VoIP.

The first pre-commercial trial network with 3G was launched by NTT DoCoMo in Japan in the Tokyo region in May 2001, using the WCDMA technology. In 2002 the first 3G networks on the rival CDMA2000 1xEV-DO technology were launched by SK Telecom and KTF in South Korea, and Monet in the US. Monet has since gone bankrupt. By the end of 2002, the second WCDMA network was launched in Japan by Vodafone KK (now Softbank). European launches of 3G were in Italy and the UK by Three/Hutchison group, on WCDMA. 2003 saw a further eight commercial launches of 3G, six more on WCDMA and two more on the EV-DO standard.

During the development of 3G systems, 2.5G systems such as CDMA2000 1x and GPRS were developed as extensions to existing 2G networks. These provide some of the features of 3G without fulfilling the promised high data rates or full range of multimedia services. CDMA2000-1X delivers theoretical maximum data speeds of up to 307 kbit/s. Just beyond these is the EDGE system which in theory covers the requirements for 3G system, but is so narrowly above these that any practical system would be sure to fall short.

The high connection speeds of 3G technology enabled a transformation in the industry: for the first time, media streaming of radio (and even television) content to 3G handsets became possible,with companies such as RealNetworks and Disney among the early pioneers in this type of offering.

In the mid-2000s (decade), an evolution of 3G technology began to be implemented, namely High-Speed Downlink Packet Access (HSDPA). It is an enhanced 3G (third generation) mobile telephony communications protocol in the High-Speed Packet Access (HSPA) family, also coined 3.5G, 3G+ or turbo 3G, which allows networks based on Universal Mobile Telecommunications System (UMTS) to have higher data transfer speeds and capacity. Current HSDPA deployments support down-link speeds of 1.8, 3.6, 7.2 and 14.0 Mbit/s.

By the end of 2007, there were 295 million subscribers on 3G networks worldwide, which reflected 9% of the total worldwide subscriber base. About two thirds of these were on the WCDMA standard and one third on the EV-DO standard. The 3G telecommunication services generated over \$120 billion of revenues during 2007 and at many markets the majority of new phones activated were 3G phones. In Japan and South Korea the market no longer supplies phones of the second generation.

Although mobile phones had the ability to access data networks such as the In-

ternet, it was not until the widespread availability of good quality 3G coverage in the mid-2000s (decade) that specialized devices appeared to access the mobile web. The first such devices, known as "dongles", plugged directly into a computer through the USB port. Another new class of device appeared subsequently, the so-called "compact wireless router" such as the Novatel MiFi, which makes 3G Internet connectivity available to multiple computers simultaneously over Wi-Fi, rather than just to a single computer via a USB plug-in. Such devices became especially popular for use with laptop computers due to the added portability they bestow. Consequently, some computer manufacturers started to embed the mobile data function directly into the laptop so a dongle or MiFi was not needed. Instead, the SIM card could be inserted directly into the device itself to access the mobile data services. Such 3G-capable laptops became commonly known as "netbooks". Other types of data-aware devices followed in the netbook's footsteps. By the beginning of 2010, E-readers, such as the Amazon Kindle and the Nook from Barnes & Noble, had already become available with embedded wireless Internet, and Apple had announced plans for embedded wireless Internet on its iPad tablet devices later that year.

In market implementation, 3G downlink data speeds defined by telecommunication service providers vary depending on the underlying technology deployed; up to 384kbit/s for WCDMA, up to 7.2Mbit/sec for HSPA and a theoretical maximum of 21.6 Mbit/s for HSPA+ (technically 3.5G, but usually clubbed under the tradename of 3G).

3G networks offer greater security than their 2G predecessors. By allowing the User Equipment (UE) to authenticate the network it is attaching to, the user can be sure the network is the intended one and not an impersonator. 3G networks use the KASUMI block cipher instead of the older A5/1 stream cipher. However, a number of serious weaknesses in the KASUMI cipher have been identified. In addition to the 3G network infrastructure security, end-to-end security is offered when application frameworks such as IMS are accessed, although this is not strictly a 3G property.

The bandwidth and location information available to 3G devices gives rise to applications not previously available to mobile phone users. Some of such applications are:

- Global Positioning System (GPS);

- Location-based services;

- Mobile TV;

- Video Conferencing

- Video on demand.

## 2.5  4G - Native IP networks

By 2009, it had become clear that, at some point, 3G networks would be overwhelmed by the growth of bandwidth-intensive applications like streaming media. Consequently, the industry began looking to data-optimized 4Th-generation technologies, with the promise of speed improvements up to 10-fold over existing 3G technologies [5]. The first two commercially available technologies billed as 4G were the WiMAX standard (offered in the U.S. by Sprint) and the LTE standard, first offered in Scandinavia by TeliaSonera.

One of the main ways in which 4G differed technologically from 3G was in its elimination of circuit switching, instead employing an all-IP network. Thus, 4G ushered in a treatment of voice calls just like any other type of streaming audio media, utilizing packet switching over Internet, LAN or WAN networks via VoIP.

A more detailed LTE overview is presented in Chapter 3.

# 3 From LTE to 5G

## 3.1 Introduction

Long Term Evolution (LTE) is designed to meet the IMT-2000 requirements set out by International Telecommunications Union - Radio communication sector (ITU-R).

LTE is a phenomenal technology: it enables operation under a vast set of conditions and still delivers excellent performance. It builds on the 3GPP GSM/UMTS cellular networks and uses Evolved-UMTS Terrestrial Radio Access Network (E-UTRAN) as its radio access. Compared to the previous 3GPP telecommunication standards, LTE marks a departure from the normal circuit switched or a combination of circuit and packet switched networks, to an all-IP/packet-based network. It is a significant advancement in cellular technologies that provides high quality experience and ensure the continuity of competitiveness of the 3G system for the future. Furthermore LTE, meets the user demand for higher data rates and quality of service and the continued requests for cost reduction in terms of both capital and operational expenditure (CAPEX and OPEX). LTE leverages on several technologies such as use of Orthogonal Frequency Division Multiplexing (OFDM) and Multiple Input Multiple Output (MIMO) antenna techniques to achieve the specified targets. Continuous improvements of these enabling technologies is the basis for the evolution of the LTE technology.

The fifth generation of mobile technology (5G) is positioned to address the demands and business contexts of 2020 and beyond. It is expected to enable a fully mobile and connected society and to empower socioeconomic transformations in countless ways many of which are unimagined today, including those for productivity, sustainability and well-being. The demands of a fully mobile and connected society are characterized by the tremendous growth in connectivity and density/volume of traffic, the required multi-layer densification in enabling this, and the broad range of use cases and business models expected. Therefore, in 5G, there is a need to push the envelope of performance to provide, where needed, for example,

much greater throughput, much lower latency, ultra-high reliability, much higher connectivity density, and higher mobility range. This enhanced performance is expected to be provided along with the capability to control a highly heterogeneous environment, and capability to, among others, ensure security and trust, identity, and privacy. While extending the performance envelope of mobile networks, 5G should include by design embedded flexibility to optimize the network usage, while accommodating a wide range of use cases, business and partnership models. The 5G architecture should include modular network functions that could be deployed and scaled on demand, to accommodate various use cases in an agile and cost efficient manner. In 5G, Next Generation Mobile Networks (NGMN) Alliance anticipates the need for new radio interface(s) driven by use of higher frequencies, specific use cases such as Internet of Things (IoT) or specific capabilities (e.g., lower latency), which goes beyond what 4G and its enhancements can support. However, 5G is not only about the development of a new radio interface. NGMN envisions 5G as an end-to-end system that includes all aspects of the network, with a design that achieves a high level of convergence and leverages today's access mechanisms (and their evolution), including fixed, and also any new ones in the future. 5G will operate in a highly heterogeneous environment characterized by the existence of multiple types of access technologies, multi-layer networks, multiple types of devices, multiple types of user interactions, etc. In such an environment, there is a fundamental need for 5G to achieve seamless and consistent user experience across time and space. Business orientation and economic incentives with foundational shift in cost, energy and operational efficiency should make 5G feasible and sustainable. 5G should also enable value creation towards customers and partners through the definition and exposure of capabilities that enhance today's overall service delivery. Enabling 5G use cases and business models require the allocation of additional spectrum for mobile broadband and needs to be supported by flexible spectrum management capabilities. NGMN and other stakeholders/partners will work together towards delivering globally and commercially available 5G solutions by 2020. This process will require a process of collaboration in the industry through existing standards development organizations (SDOs), or potentially new collaboration forms like open source.

A brief survey on the technologies that made LTE a tipping point for the mobile communications and the key-enable is presented in the Section 3.2.

Section 3.3 core and the radio access network solutions for the Next Generation 5G system are presented. Moreover the main architecture innovations (i.e. Mobile-

access Edge Computing – MEC) and strengths (i.e. virtualization, reliability and slicing) that will characterize the Next Generation 5G system are introduced.

## 3.2 LTE network architecture

LTE network architecture is a generally simplified access network which marks a total departure from previous standards, characterized by the absence of a circuit-switched domain. It employs a non-hierarchical (distributed) structure. The LTE network architecture incorporates new network elements [6].

As shown in Fig. 3.1, LTE network architecture can be sub-divided into three major groups: air interface (orange areas in Fig. 3.1), radio access network (green area in Fig. 3.1) and core network (blue area in Fig. 3.1). Transmission of data and control information between the user equipment (UE) and the evolved base stations (eNBs) take place within the air interface. LTE uses various mechanisms within the air interface to provide highly reliable and efficient means of carrying out these operations.

The RAN of LTE consists only of a network of fully interconnected eNBs; hence the network is described as being flat or distributed. This RAN is called the E-UTRAN i.e. the Evolved-UMTS Terrestrial Radio Access Network. It is an evolved RAN from UTRAN, used by 3G networks but in LTE, all Radio Network Controller (RNC) functions are transferred to the eNBs. Some of the functions of the eNB include:

- Radio Resource Management: This involves functions such as scheduling, dynamic allocation of resources, radio bearer control and mobility control;

- IP Header Compression;

- Security;

- Connection of users to the core network.

The core network of LTE differed significantly from previous standards. All others had their core networks either entirely circuit switched or split into circuit switched domain and packet switched domain, but LTE core network is entirely packet switched and it is called Evolved Packet Core (EPC). The EPC in conjunction with the E-UTRAN is called the Enhanced Packet System (EPS), whose details have been defined by 3GPP's study of System Architecture Evolution (SAE).

A Summary of the functional elements of the EPC are outlined below [7]:

**Figure 3.1:** LTE network architecture

- Mobility Management Entity (MME): this handles user authentication, it tracks and maintains the location of a user equipment, performs signalling operations, MME selection for inter-MME handovers.

- Serving Gateway (S-GW): while the MME handles control distribution functions, the S-GW handles data bearer functions where it handles user data functionality, routes and forwards data packets to the P-GW, performs mobility anchoring for inter-3GPP mobility and is responsible for lawful interceptions.

- Packet Data Network Gateway (P-GW or PDN-GW): It handles packet filtering for every user, allocation of IP addresses to the UEs, supports service level charging by collecting and forwarding call data records, handles DL data rate enforcement to ensure that a user does not surpass his traffic rate subscription level, provides interworking for the user plane, between some 3GPP access systems and all non-3GPP access systems, supports QoS differentiation between multiple IP flows. It is also capable of handling multiple lawful interceptions of user traffic to promote government intelligence services fighting criminal activities. The P-GW enforces PCRF policies.

- Home Subscriber Server (HSS): this is a major database, which houses all subscription-related information, to perform call control activities and ses-

sion management functions.

- Policy and Charging Control Function (PCRF): The PCRF ensures QoS regulation within the network based on definite policies. It is responsible for framing policy rules from the technical details of Service Date Flows (SDFs) that will apply to the users' services, and then forwarding these rules to the P-GW for enforcement.

- Evolved Packet Data Gateway (ePDG): The ePDG provides interworking with un-trusted non-3GPP IP access systems. It ensures security by having a secured tunnel between the UE and the ePDG. It can also function as a local mobility anchor within un-trusted non-3GPP access networks.

As observed in Fig. 3.1, LTE uses interfaces as indicated for communication between its entities. In general, LTE network architecture implements a simplified, flat all-IP architecture which leads to reduced latency, reduced CAPEX and OPEX, increased scalability and efficiency among other benefits.

## 3.2.1 Multiplexing/Multiple Access Mechanism in LTE

The aim of multiplexing/multiple access mechanism is to share scarce resources in order to achieve high capacity.

*Multiplexing* is a mechanism by which multiple signals are transmitted at the same time in form of a single complex signal over a shared medium and then recovering the individual signals at the receiving end.

*Multiple access* mechanism define instead how the channel is shared in a finite frequency bandwidth i.e. it controls how to share the radio resources efficiently. These operations take place within the radio air interface of the LTE network.

Majority of the striking features of LTE is made possible by its use of the Orthogonal Frequency-Division Multiplexing (OFDM) data transmission multi-carrier modulation technique [8]. The OFDM divides a high bit-rate data signal into several parallel low bit-rate data signals which are then modulated using an appropriate modulation scheme. The "low bit-rate multi-carrier" technique of the OFDM, with a cyclic prefix added to it, makes the transmission robust to time dispersion on the radio channel without the need for advanced and complex receiver channel equalization. A cost and a power reductions of terminal equipment are then obtained. OFDM is also used due its resilience to multipath delays and spread,

its capability for carrying high data rates and its ability to support both FDD and TDD schemes.

A derivative of OFDM, the Orthogonal Frequency-Division Multiple Access (OFDMA) that combines functionalities of FDMA and TDMA is used in the LTE downlink. With OFDMA, the User Equipment (UE) that acts as receiver in the downlink and therefore it does not have multiple access problems in terms of collision, gets scheduled to a time slot and a frequency group to send information. By means of such scheduling the system is made resilient to frequency-selective fading. Using OFDMA, LTE can use channel-dependent scheduling to take advantage of the channel variations resulting in a more efficient use of available radio resources.

In the uplink, the UEs transmit to the Base Station (BS). Due to the high peak-to-average ratio (PAR) of OFDM i.e. the high amount of power required by the RF power amplifier to push out the RF signal from the UE antenna to the BS, 3GPP was forced to adopt a different transmission scheme for such link. Single-carrier FDMA (SC-FDMA), a hybrid scheme that combines the low PAR feature of single-carrier schemes with the resilience of multipath interference and the flexible subcarrier frequency allocation of OFDM technology, was the solution. The low PAR characterizing the SC-FDMA allows high RF power amplifier efficiency in the UEs that leads to a UE battery consumption reduction [9].

### 3.2.2 Coding and Modulation in LTE

The reduced latency and high throughput of LTE is traceable to a number of mechanism implemented in it. The physical/MAC layer of LTE adopts two key techniques: Hybrid Automatic Repeat reQuest (HARQ) and Adaptive Modulation and Coding (AMC). These two techniques work together to give a very adaptive transport mechanism in LTE [10].

HARQ is a technique for both error detection and correction by identifying when transmission errors occur and facilitating re-transmission from the source thereby ensuring the reliably data transportation from one network node to another. Please be advised that LTE uses Type-II HARQ protocols. To handle re-transmission errors, LTE uses two loops. A fast HARQ inner loop to take care of most errors a robust selective-repeat ARQ outer loop to take care of residual errors.

LTE demonstrates a dynamic resource allocation through link adaptation. Link

adaptation is achieved using the AMC mechanism, with the aim of improving data throughput in a fading channel. AMC works by varying the downlink modulation technique depending on the channel conditions of each user. Given a good channel condition, the LTE system can use a higher order modulation scheme (64-QAM with 6 bits per symbol) or reduced channel coding, making the channel more spectrally efficient and resulting in higher bit rates. As the channel become noisy due to signal fading or interference, the system selects a lower modulation technique (QPSK or 16-QAM with fewer bits per symbol) or stronger channel coding to obtain a signal more robust at the expense of the bit rate.

### 3.2.3 Radio Channel Bandwidth in LTE

LTE is not only able to operate in different frequency bands, but can be also implemented using different spectrum sizes. This feature make possible to harness the global wireless market and align with regional spectrum regulations.

LTE implements a scalable radio channel bandwidth from 1.4 to 20 MHz with a subcarrier spacing of 15 KHz. The 20 MHz bandwidth is required for optimum performance and to cope with the growth of the mobile internet. 3GPP has specified the LTE air interface to be "bandwidth agnostic" to allow the physical layer to adapt different spectrum allocation without severe impact on system operation.

### 3.2.4 LTE Technological Advancements

There are two groups of technological advancements on LTE Release 8, namely LTE Release 9 and LTE Release 10 [11].

Transmission Mode 8 (TM8) and Dual Layer Beamforming were added in LTE Release 9, that also focuses on following features enhancing the LTE Release 8 core network.

- Location, broadcast and IP Multimedia Subsystem (IMS) emergency services using GPRS and EPS;

- Support of circuit switching services over the LTE EPS;

- eNB considerations focusing on security, QoS, charging and access restrictions;

- IMS evolution.

The LTE Release 10 is an evolution of LTE to meet the IMT-A requirements defined by ITU. It is known as LTE-Advanced (LTE-A) and it focuses on higher capacity as following:

- Downlink peak data rate of 3 Gb/s and uplink peak data rate of 1.5 Gb/s;

- Higher spectral efficiency on the downlink (e.g. 30 bps/Hz);

- Increased number of simultaneously active subscribers;

- Improved performance at cell edges.

To achieve the above-mentioned feats, the three following new techniques are considered in LTE-A [12]:

- Carrier Aggregation (CA): The most basic way to increase capacity is add more bandwidth. LTE-A is increased in bandwidth by means of aggregation of up to five component carriers of different bandwidths to form a maximum bandwidth of 100 MHz. CA can be used in both FDD and TDD schemes.

- Enhanced multiple antennas techniques: in LTE-A a ninth transmission mode called Eight Layer Spatial Multiplexing (8 x 8 MIMO) and a second transmission mode (4 x 4 MIMO) are added to the downlink and uplink respectively.

- Relay Nodes (RN): Relay Nodes bring the possibility of efficient heterogeneous network in LTE-A. The Relay Nodes are low power base stations that provide enhanced coverage and capacity at cell edges. By means of RN, connectivity without the need of optical fibres can be provided to remote areas.

- Coordinated Multipoint (CoMP) Transmission/Reception: This feature was finalized in Release 11. In this technique, multiple transmit and receive points provide coordinated transmission/reception. This transmission/reception is carried put jointly and dynamically across multiple cell sites, same site or within same or different eNBs. The primary purpose of CoMP is to improve the performance at cell edge.

The on-going enhancements in Release 13 include additional features for allowing the LTE to operate in the unlicensed spectrum and the expansion of the CA framework to support more than five component carriers. Other enhancements in Release 13 include:

- Enhancements for Machine-Type Communications (MTC) defining a new low complexity UE type that supports reduced support for downlink transmission modes, reduced bandwidth, reduced transmit power and very long battery life to support the rise of Internet of Things (IoT) markets;

- Improving multi-user transmission techniques using superposition coding for increasing the LTE system spectral efficiency;

- Use of full dimension MIMO/Elevation Beamforming for improved spectral efficiency by the use of higher dimension MIMO of up to 64 antennas at the eNB and utilizing the vertical dimension for MIMO and beam-forming operations.

- Improved indoor positioning accuracy and support for Single-cell Point-to-Multipoint (SC-PTM).

Table 3.1 gives a summary of the key characteristics of LTE at its inception and the current features of LTE as today, LTE-A.

| Parameter | LTE | LTE-A |
|:---:|:---:|:---:|
| Frequency Band | Country-dependent | Country-dependent |
| Downlink Peak Data Rate | 100 - 326 Mbps | 1 - 3 Gbps |
| Uplink Peak Data Rate | 50 - 86 Mbps | 500 Mbps - 1.5 Gbps |
| Channel Bandwidth (MHz) | 1.4, 3, 5, 10, 15, 20 | UP to 10 MHz |
| Peak Spectral Efficiency | 16 bps/Hz (DL) | 30 bps/Hz (UL) |
| Latency | 10 ms | Less than 5 ms |
| Duplex Method | FDD and TDD | FDD and TDD |
| Multiplexing | OFDM | OFDM |
| Multiple Access Method | OFDMA in DL SC-FDMA in UL | OFDMA in DL SC-FDMA in UL |
| Modulation Scheme | QPSK 16-QAM 64-QAM | QPSK 16-QAM 64-QAM |
| Multiple Antenna Technique | Up to 4 x 4 MIMO (DL) | 8 x 8 MIMO in DL 4 x 4 MIMO in UL |

**Table 3.1:** LTE vs LTE-A

## 3.2.5 Towards the 5G

The most obvious paths of evolution towards 5G radio access are improved spectrum efficiency, network densification and spectrum extension. As stated in [13], currently deployed networks use 1-3 GHz frequency band which eventually fall short of meeting the multi-gigabit requirements of future communication services such as Ultra-High Definition Video (UHDV) [14]. The millimeter wave (mmWave) frequency band (from 30 to 300 GHz) offers huge bandwidth and consequently spectrum extension for mobile networks. mmWave communications particularly in the 28, 30, 60 GHz and E-band (71-76 and 81-86 GHz) bands will play a critical role in 5G applications such as small cell access, cellular access and wireless backhaul. Some of the key radio access technologies that will pave the way for 5G mobile communications include:

- Further improvements to low power small cells to provide network densification;

- The use of massive MIMO and large number of miniaturized antennas at mmWave frequencies to provide significant increase in spectrum efficiency and user throughput;

- Use of access techniques such as Filtered OFDM and Sparse Code Multiple Access (SCMA) to improve system efficiency, support energy saving, reduced latency and massive connectivity;

- Use of more efficient coding schemes such as Polar codes, which can achieve Shannon capacity using simple encoder.

- Use network coding for interference management and for the security, throughput and robustness for routing information through the network improvement;

- Use of Full Duplex to support bi-directional communications without the use of time or frequency duplex, to double the system capacity and reduce the latency;

- Use of self-organizing network operation for a cost effective management of the massive network densification.

The LTE technology has undergone a significant evolution from its first release, which was aimed at meeting the IMT-2000 requirements to achieving and even

exceeding the IMT-Advanced (4G) requirements. These technologies played a critical role in the frequency spectrum below 6 GHz that has been allocated for mobile communication at the World Radio Conference (WRC) in 2015. However, for the spectrum band above 6 GHz which is expected to be allocated at WRC in 2019, a new radio access technology is necessary. Thus 5G is the next frontier of a broader ICT ecosystem that will enhance mobile internet and empower Internet-of-Things will be heterogeneous across frequency spectrum.

## 3.3 The 5G Next Generation Mobile Network

It is expected that mobile and wireless traffic volume will increase a thousand-fold over the next decade which will be driven by the expected 50 billion connected devices connected to the cloud by 2020 and all need to access and share data anywhere and anytime.

With the rapid increase in the number of connected devices, some challenge appear which will be responded by increasing capacity and improving energy efficiency, cost and spectrum utilization as well as providing better scalability for handling the increasing number of the connected devices. For the vision of all-communicating world relative to today's network, the overall technical aim is to provide system idea that supports [15]:

- 1000 times increased data volume per area;

- 10 to 100 times increased number of connected devices;

- 10 to 100 times increased typical data user data rate;

- 10 times extend battery life;

- 5 times reduced end to end latency.

To meet such demands of the user and to overcome the challenge that has been put forward in the 5G system, a drastic change in strategy of designing the 5G wireless cellular architecture is needed. New novel concepts such as Software defined networking (SDN), Network functions virtualisation (NFV), Network slicing and functional split will have to be considered in the 5G system deployment.

All these novel aspects will be introduced in the next Chapters.

**Figure 3.2:** Services in Future 5G Networks

### 3.3.1 5G Core Network Architecture

To use 5G to support a wide range of services with different demanding performance requirements, and to insert 5G into new industrial value chains, 5G requires a new system architecture and in particular a new core network known as Next-Generation Core (NG-Core). A "cloud native" core will enable operators to achieve flexibility, scalability, reliability and performance needed to meet 5G service targets [16].

The NG-Core should support and enable the wide range of services envisioned for 5G across several market segment, as outlined in Fig. 3.2.

There are several principles that illustrate what 5G NG-Core will look like and can inform specification and development. These are the follows:

- **Meet Demanding KPIs**. Performance is critical to ensure availability, latency, reliability, user experienced data rates and area traffic capacity.

- **Access Independent**. The new core should support multiple access technologies and key functions should be decoupled from access when possible and appropriate.

- **Flexible, Scalable & Programmable**. To adapt to change and support dynamic services using network slicing, the new core will sue cloud-native architectures and software technologies.

- **Support Real-Time & Non-Real-Time Services**. NG-Core should support highly-dynamic and variable services with appropriate quality of experience (QoE).

- **Interwork with Existing Networks**. This includes existing 4G core networks, multiple non-3GPP access types, ad service-layer technologies such as IMS/VoLTE and IoT platforms.

Anyway, in the end, 5G needs a new, software-centric architecture designed to operate in a modern, cloud-native networking environment. This is true not only in principle, but also in terms of specific capabilities and services provided by the NG-Core. These include:

- **Access Agnostic Core**. The existing core is not access-independent. While non-3GPP access is supported by the LTE EPC, it requires integration of specific equipment to connect, for example, Wi-Fi into the mobile core. In NG-Core, mobility management may only be instantiated as needed; fixed access would only need a subset of the NG-Core features to operate. The addition of unlicensed radio will also generate new interworking requirements.

- **Connectionless Services**. 5G will move from the "always-on" model in 4G to an "always-available" model in which connectivity management is used as needed, for example, for session continuity or mobility. This is particularly important for IoT devices that, to save battery life, unscheduled uplink transmissions on both the data and control plane to transmit data and signaling traffic.

- **Ultra-Low Latency & Mission-Critical Services**. Some advanced 5G services, such as tactile Internet and industrial control systems, require ultra-low latency. This will generate a new architecture that limits the physical distance between access and core and makes use of distributed forwarding elements. A new session and service continuity model, as well as a new QoS model, that efficiently enables guaranteed services will also be needed.

The Cost and performance of a Web-scale infrastructure has inspired operators to apply cloud principles to their own networks: commercial-off-the-shelf (COST) hardware, distributed processing, centralized control, model-driven configuration, automation, etc., now inform the evolution of operator networks. Table 3.2 compares some high-level differences between virtualized networks and cloud-native

native networks. Given the timeline and the nature of 5G services, NG-Core will be designated for, and deployed on, cloud infrastructure.

| Virtualized | Cloud-Native |
|---|---|
| Manual Operations; limited, workflow/script-driven automation | Extreme automation, model and policy-drive, "post DevOps" |
| Investment driven by standard refresh cycles | New Investment driven by new business/models (e.g. digital, 5G, IoT) |
| COTS + hypervisor + VNF | Multi-vendor, horizontal, interoperable cloud components |
| COTS servers, Linux, hypervisors, OVS/VPP, Openstack | PaaS, machine learning, hyper-converged servers, "compact" DCs |
| Handful of VNs (limited scaling) | Dozens to hundreds of VNS (dynamic, ephemeral networks) |
| Software designs and reliability/redundancy schemes based in physical resources | Software designs built with cloud principles and cloud reliability and recovery mechanisms |

**Table 3.2:** Virtualized vs. Cloud-native Networks

The 5G functional architecture, including the NG core itself, is now under development in 3GPP Release 15. This specification work is fundamental to how the 5G NG-Core will be designed and then deployed and operated in commercial networks. The specification work covers both non-standalone (NSA) mode, in which the 5G radio is integrated with LTE, and standalone (SA) mode, which enables 5G radio to be deployed using NG core without any LTE dependencies.

The two proposed architecture are shown in Fig. 3.3 and 3.4 respectively. The first one is a point-to-point architecture, (see Fig. 3.3), that can be thought of as a traditional 3GPP architecture. It splits control- and user-plane and separates access and mobility management (AMF) and session management (SMF) to enable independent evolution and scaling. On the access side, it minimizes core network dependencies by specifying common interfaces for 3GPP and non-3GPP access types. In the user-plane, it enables advanced features, such as User plane Function (UPF) branching, if the service requires. In short, it is familiar, yet also novel enough, to meet the needs of 5G in the near-term, and is flexible enough to evolve over time. One challenge with this model is that a change in topology occurs when a new function is added, and this requires new interfaces to be es-

**Figure 3.3:** Point-to-Point NG-Core Architecture

tablished between neighboring functions. Since the service consumption is tied to the network function, this is one of the main factors that makes mobile core networks relatively difficult to change and adapt. NG core should be virtualized and cloud-based to make changes of this kind faster and easier. Because the Point-to-Point NG core architecture is essentially nodal, it is easily imaginable that the functional elements belonging to it will be developed as discrete pieces of equipment or as discrete VNFs. A truly cloud-native model would probably lean toward a network function as-a-service model.

The second proposed architecture defined as a Service Based Architecture (SBA), shown in Fig. 3.4, is, arguably, more aligned with modern cloud principles. The functional elements are the same of the Point-to-Point NG-Core architecture and also the N1, N2, N3, N4 interfaces to the control-plane are unchanged. The difference is that rather then having predefined interfaces between the control-plane nodes, the functions themselves present "service interfaces" (APIs) to each other on an on-demands basis. The NF Repository Function (NRF) provides registration and discovery mechanisms to enable the different control-plane components to communicate directly. The service-based NG-Core architecture reflects the idea of a *Network Cloud OS*, where network services are composed using a library of functions hosted in the cloud and chained together to create the end-to-end service. The challenge is actually implementing this model, as it relies on the availability and maturity of cloud networking technologies outside the 3GPP domain. Indeed, for example, resource and service orchestrators needed to enable this are part of the NFV cloud platform environment that NG-Core will be deployed in.

These two options are not directly competitive, but instead can be thought of as

**Figure 3.4:** Service-Based NG-Core Architecture

two different ways of representing a common set of functional elements. It is possible, perhaps likely even, that operators will start with a point-to-point architecture and then migrate to a service-based architecture over time.

The main functional elements of the above proposed NG architectures are briefly described in the following:

- AUthentication Function (AUSF): keeps a key for reuse, derived after authentication, in case of simultaneous registration of a UE in different access network technologies, i.e. 3GPP access networks and non-3GPP access networks such as IEEE 802.11 Wireless Local Area Network (WLAN).

- Access and Mobility Management Function (AMF): a control-plane component that manages access control and mobility. It handles all the 5G signaling coming from and going to the UE, supports user access to the network and manages mobility by interacting with the UE and with other NFs (e.g., SMF, AUSF, etc.). The AMF will likely also include network slice selection functionality. The AMF contains part of the LTE MME functionality.

- Session Management Function (SMF): is the control part of a PDU session. That is, it configures NG tunnels, allocates IP addresses with DHCP, and configures traffic steering (e.g., towards a third party or an edge cloud). The SMF contains parts of the LTE MME and P-GW.

- User Plane Function (UPF): it handles the NG user plane (NG-U) tunnel forwarding and the related data path services, such as anchoring for handover, QoS, and traffic policy enforcement. The UPF contains parts of the LTE SGW and PGW functionalities. Multiple different UPFs, in distributed and centralized locations, can be used by operators, according to the services type.

- Policy Control Function (PCF): provides a common policy framework incorporating network slicing, roaming and mobility management to other control plane functions, such as SMF.

- Unified Data Management (UDM): component used for storage the subscribers data and profiles (i.e. credentials, identifiers, AMF details, and SMF assignments) for the current session. The underlying idea of the UDM is to create a central database for UE configuration information, so that the NFs can be designed as stateless services, improving architectural agility The UDM contains part of the LTE HSS functionality.

- NR Repository Function (NRF): this is a new functionality without equivalent in LTE and it is present only in the service-based architecture. It provides registration and discovery functionality so that network functions (NFs) can discover each other and communicate directly. When it receives an NF discovery request from a NF instance, it provides the discovered NF instances.

- Application function (AF): resembles an application server that can interact with the other control-plane NFs. AFs can exist for different application services, and can be owned by the network operator or by trusted third parties. For instance, the AF of an over-the-top application provider can influence routing, steering its traffic towards its external edge servers.

- Network Exposure Function (NEF): exposes the capabilities of networks and network/UE events for third-party, application function, edge computing, and other purposes. It is not present in LTE and it is implemented only in the service-based architecture.

- Data Network (DN): stands for e.g. operator services, Internet access or 3rd party services.

## 3.3.2 5G Radio Access Network

The 3GPP 5G RAN architecture, specified in Release 15 and known as NG-RAN, introduces new terminology, interfaces and functional modules. The NG-RAN consists of a set of radio base stations, called gigabit NodeB (gNB) connected to the 5G NG-Core and to each other. The gNB incorporates three main functional modules [17]: (i) the Centralized Unit (CU), (ii) the Distributed Unit (DU), and the (ii) Radio Unit (RU), which can be deployed in multiple combinations. In particular in the CU, the gNB protocol stack upper layers (e.g. Packet Data Convergence Protocol - PDCP, Radio Resource Control - RRC) are hosted; in the DU, the lower layers (e.g. Physical - PHY, Medium Access Control - MAC, Radio Link Control - RLC) are hosted, and finally the RU hosts the Radio Frequency (RF) functionalities. RU and DU communicate using a fronthaul interface (also called Fronthaul I - F1) while DU and CU communicate through a midhaul interface (also called fronthaul II). The F1 interface is expected to be interoperable across vendors. The CU can be further disaggregated into the CU user plane (CU-UP) and CU control plane (CU-CP), both of which connected to the DU over the F1-U and F1-C interfaces respectively. This new 5G RAN architecture is described in 3GPP TS 38.401 while several functional split options and their specific requirements in terms of data rate and latency have been planned by 3GPP technical report TR 38.801.

NG-RAN is therefore a logical architecture that can be implemented and deployed in different ways, according to an operator's requirements and preferences. As shown in Fig. 3.5, the base station can be deployed as monolithic unit deployed at cell site, as in a classic cellular networks, or split between the CU, DU and RU. The CU-DU interface is a Higher-Layer Split (HLS), which is more tolerant to delay. The DU-RU interface, which is not yet standardized, is a Lower-Layer Split (LLS), which is more latency-sensitive and demanding on bandwidth, but may offer improved radio performance across a coverage area due to coordination gain. CUs, DU, and RUs can be deployed at locations such as cell sites (including towers, rooftops and associated cabinets and shelters), transport aggregations sites and "edge sites" (e.g., central offices or local exchange sites).

The choice between LLS and HLS is a critical decision in NG-RAN architecture. The challenge is that the trade-offs are not always clear-cut, and it can be hard to determine which model will be economically optimal for each carrier. Moreover, it may make sense to use different models for different regions, for example rural coverage versus urban capacity. There will be also variation for different use cases.

**Figure 3.5:** Flexibility for 5G RAN Functional Units



**Figure 3.6:** Flexibility for 5G RAN Functional Units

Indeed an ultra-low-latency factory automation network or an Advanced Driving Assisted use case would require DU/CU close to, or integrated with, the RU. Fig. 3.6 shows the high-level trade off between coordination gain from centralization and the latency and bandwidth requirements in the transport network.

Usage patterns, device capabilities, operating costs, RF strategy and the existing network footprint and capabilities are all influential on the decision to determine the optimal RAN topology and associated operating model. A lower layer split between the DU and RU is demanding in terms of bandwidth, latency and packet loss, and to deploy a DU at a distance from the cell site typically requires the operator to invest in fronthaul fiber infrastructure. A higher-level RAN split, with DUs deployed at the cell site is more tolerant in terms of transport performance, but perhaps offers less coordination gain and, potentially, higher OPEX in the longer run.

In Fig. 3.7 are depicted the primary functional elements of the 5G RAN (RU, DU and CU) and how they are mapped to the transport network. The 5G RAN deployment architecture is therefore codependent on RAN design, the transport network and end-users services. Please note that the split was also possible in 4G, but in 5G it is part of the architecture that can support a number of deployment

**Figure 3.7:** Flexibility for 5G RAN Functional Units

options (e.g. co-located CU-DU deployment is also possible).

### 3.3.3 Virtualization in 5G

The 5G network will require virtualization both at the network's core and at the network's edge. Therefore, to address the strong technological advances that the forthcoming 5G standard promises, the mobile operators are already seeking to enhance their investments to handle ever-increasing consumer and business connectivity demands, so it is only natural that they are turning to infrastructure-enhancing technologies such as network virtualization (NV), software-defined networking (SDN), and network functions virtualization (NFV) to evolve their business.

Therefore, the future 5G network architecture is expected to be massively based on Network Functions Virtualization (NFV).

In the network core, 5G virtualization will use network slicing to support multiple virtual networks over one physical network infrastructure. Network slicing permits the logical separation of a network so that each slice provides unique connectivity, but all slices run on the same shared infrastructure. In this way, 5G virtualization provides a new level of flexibility, allowing operators to, for example, devote a network slice to certain kinds of devices. To efficiently support certain sets of services, each network slice will be able to access different types of resources, such as infrastructure (e.g., VPNs, cloud services) and virtualized network functions (VNFs). With 5G virtualization, operators will be able to create

custom networks with unique sets of capabilities.

Virtualization will be an essential component also at the network edge, namely, the virtual partitioning of the mobile RAN. Operators will be able to create unique services that are customized for various use cases such as IoT, automated cars, streaming video, remote health care, and so on. They can create virtual networks for those applications that boast separate blends of performance, capacity, latency, security, reliability, and coverage. The virtualization of NG-RAN components (i.e., CU and DU) allows to move toward a Virtualized RAN (VRAN) and achieve the full potential of cost saving with rapid deployment of new services. Because of the additional hypervisor layer (the fundamental building block of virtualization), the midhaul segment requirements in terms of latency and jitter, reported in 3GPP TR 38.801, may change and become more stringent. So far, several works studied the impact of virtualization on physical infrastructure sharing, isolation, cost, and energy saving of Long Term Evolution (LTE) networks. However, the estimation of the effect of virtualizing the NG-RAN components on the fronthaul/midhaul latency and jitter (i.e., the maximum allowable latency) has not been conducted in details yet. Chapter 5 is completely dedicated to the estimation of the effect of virtualizing the NG-RAN by means of experimental performance measurements.

## 3.3.4 Multi-Access Edge Computing achieves 5G goals

Both Multi-access Edge Computing (MEC) and 5G are considered disrupting technologies on their own but combined they will become a powerful force in the world of computing. MEC uses the edge of the network to bring computing closer to the data center, which reduces latency and increases connection speeds. 5G promises to have estimated network speeds as fast as 10 Gb/s. The emergence of 5G networking capabilities will increase the number of connected devices on a network, which spurs the need for edge computing to help distribute networking demands.

Applications that rely heavily on a consistent network connection, rapid deployment, and low latency include burgeoning technologies such as artificial intelligence (AI), Internet of Things (IoT), and virtual reality (VR). MEC and 5G networking together allow for the simultaneous usage of a massive number of connected technologies without incurring network outages due to traffic bottlenecks.

In other words, MEC achieves 5G goals. In MEC, computing capabilities are pushed closer to the radio access network (and, in turn, closer to subscribers),

enabling low-latency and high-bandwidth access to content, applications and services. The distributed MEC architecture also makes it ideal for supporting high volumes of connected devices, which will generate even higher volumes of data interactions. Backed by industry leaders participating in ETSI's Mobile Edge Computing industry standardization group (ISG), MEC provides a standard-based approach to making significant progress towards the 5G in LTE networks today.

MEC can be leveraged today to create an open ecosystem and growth engine inside the NG-RAN in close proximity to mobile users. Leveraging a Mobile Edge Computing platform enables mobile operators to quickly and effectively deploy new revenue generating services for content delivery, Internet-of-things (IoT) connectivity, retail, and enterprise applications.

Some advantages that MEC and 5G together will produce are following:

1. The utilization of microdata centers that deliver continuous edge networking for mobile computing.

2. The localization of the edge computing that significantly decrease end-to-end latency when used in conjunction with 5G.

3. Enhanced quality of service to end users, particularly when it comes to video streaming and IoT technology such as connected cars, enabled through the use of network slicing. Network slicing is a 5G innovation that allows a network's functions to be divvied up depending on an application's demand on the network. For instance, a connected car requires reliable, speedy connection for its optimal performance. The network may be sliced to guaranteed that the connected car receives it's high networking requirements while reducing the size of the networking slices for other applications that are not as demanding on the network.

4. The MEC server platform will be able to deliver real-time analytics and machine intelligence.

5. Data analytics conducted at the end of a network cuts down the amount of data sent upstream, resulting in decreased cost to the enterprise and/or service provider.

In Chapter 6 the benefits of MEC deployment in the Next Generation 5G networks and in particular in the automotive scenario (i.e. connected vehicles, vulnerable road user safety and vehicles infotainment), are deeply described and analysed. Two deployed MEC-based applications are also presented.

## 3.3.5 Reliability in 5G

One of the design requirements of the Next Generation 5G system is guaranteeing nearly 100 percent up-time and ultra-consistent network services.

5G key performance indicators (KPI) specified Ultra Reliable Low Latency Communication (URLLC) service requirements of no more than $10^{-5}(0.001/\%)$ of 20 byte packets can fail to be delivered by 1ms. This is a packet and latency service level commonly provided by network service providers. This is about 99.999 percent reliability and 99.999 percent availability. The reliability performance indicator fails if:

- Too many packets are lost

- Too many arrive too late

- The packets have errors

This low-packet error rate is helpful, but networks must be available 24 by 7, every day without fail. Customers require - guaranteed service - not only for latency and reliability, but also connectivity.

The most vulnerable segment between client devices or the network edge and the carrier network is the last few hundred feet of communications running over wireless technologies. Everyone has experienced LTE dropped calls: the 4G specification was designed for smartphones and not for critical network services. Wireless links notoriously fluctuate in quality, leading to variable packet error rates, which are orders of magnitude higher than wired (i.e., "Can you hear me now?").

5G has increased the vulnerability of the wireless technologies with the inclusion of high frequencies from the mmWave spectrum. mmWaves are line-of-sight directional and easily blocked. Rain and other obstructions routinely cause network packet errors, often resulting dropped links and outages. These network challenges can be managed with new 5G technologies like beam-forming, network fast-failover, SDN technologies and wireless packet retransmission protocol features.

Service confidence levels can be achieved with a well-planned network architecture. These include radio access technology (RAT), measuring the spatial availability of a wireless link for a service-relevant confidence level, temporal availability and falling back to other RATs (e.g., LTE).

5G wireless networks are implemented using a combination of diversity coding and network coding schemes. These represent a forwarding architecture for

**Figure 3.8:** Flexibility for 5G RAN Functional Units

wireless mesh networks, which improves throughput and consistent network performance. A coding layer is inserted between the IP and MAC layers. This provides robust link-failure recovery and near-instantaneous packet recovery. Traffic can be routed across diverse paths for greater network throughput and reduced re-transmissions using error-control techniques.

3GPP Release 15 includes a few reliability improvements over LTE. The flexible frame structure with various options for sub-carrier spacing, modulation and coding provide a solid foundation at the physical layer. Hybrid automatic repeat request (HARQ) and ARQ are included and will initiate "retries" on corrupted data (see Fig. 3.8). In addition, the 3GPP specification includes reference signals to improve synchronization, enhancing demodulation efficiency and significantly reducing packet data corruption and errors.

5G network architectures have redundant hardware, software and network links with automated failover technologies. This guarantees that no system or network outage will interrupt network services.

In Fig. 3.9, multiple Central Units (CUs) have multiple redundant network connections. Each CU can take over the functions of a failed CU. The Remote Units (RUs) wireless towers are midhaul-connected to multiple DUs stations with dark fiber networks, and each RU tower has wireless network connections to other RU towers. Each RU tower covers an area that overlaps with other towers. A loss of one tower will not leave a gap in service coverage.

Reliability to business or organizational facilities can be provided by redundant and diverse connections to the 5G carrier networks. Wireless connections should

**Figure 3.9:** Highly-available 5G Network Topologies

be made to two or more 5G base stations. Failed network links or components can be rerouted to back-up or redundant network links.

High reliability is one of the key design specifications for 5G. Mission-critical services such as connected robotic factories, remote surgery, patient care and driverless cars are required to be constantly connected to the network. As such, 5G network outages could be life threatening. The 5G architecture has been designed to be highly available with redundancy integrated into every component. Network architects can extend this level of reliability and availability with a well-planned architecture that includes redundant 5G network access paths and hardware.

Chapter 7 deals with the issue of reliability in 5G system analysing both core and RAN reliability and proposing different resilience schemes and protection and restoration algorithms.

### 3.3.6 5G Slicing

In 5G, the term slicing refers, in general, to the possibility for different customers (usually called tenant) to share the same physical network.

Thanks to the softwarization of networks according to the Network Function Virtualization (NFV) concept and the programmability of network connectivity through Software Defined Networking (SDN), new network and service capabilities can be envisioned by integrating networking, computing and storage resources while serving a multitude of tenants. Each tenant is assigned a logical network that can satisfy its requirements. Survivability is one of the most important requirements especially for vertical applications requesting Ultra Reliable Low Latency Communications (URLLC).

With the advent of NFV and SDN a novel network scenario is envisioned enabled by network deployments into the cloud also extended to the network edge and by programmability of network connectivity through network controllers. This trend known as softwarization is enabling new unique network and service capabilities by integrating networking, computing and storage resources into one programmable and unified infrastructure while serving a multitude of distributed smart devices and applications (e.g., robots, drones, smart vehicles). As result, current communications network scenario is moving from having a separate network for each application (e.g., fixed telephone network, mobile telephone networks, Internet access) to a single network shared by different applications or verticals.

Network Slicing is a key feature of the 5G System that allows Operators to flexibly structure the network resources to match the services offered to subscribers, third-party customers, including the roaming scenario.



**Figure 3.10:** Slice Concept

The concept of slicing emerged as a way of setting up several logical networks for different verticals on the same physical network. Each vertical is then assigned to the logical network that guarantees the required QoS. Such setup potentially allows communication providers to save capital and operating expenditures (CAPEX

and OPEX). However, as for any shared medium, guaranteeing the required QoS to network slices sharing the same physical network is not a trivial task and remains an open issue. In particular, slice control and management planes shall be designed for slice provisioning and dynamic reconfiguration and the data plane shall guarantee each slice requirements (e.g., QoS requirements, slice isolation, etc.).

Several Standard Developing Organizations (SDOs) are focusing on the network slicing concept [18]. The Next Generation Mobile Networks (NGMN) alliance defines a Network Slice Instance (NSI) as a set of network functions and resources forming a complete instantiated logical network to meet certain network characteristics required by the Service Instance. In [19] the network slicing concept consists of three layers depicted in Fig. 3.10: Service Instance Layer, Network Slice Instance Layer, and Resource layer. The Service Instance Layer represents the services (i.e., end-user or business services) which must be supported. The Network Slice Instance Layer provides the network slice instances with specific network characteristics that are required by the related Service Instances (e.g., Enhanced MBB, M2M, Enterprise and Industry). The Resource Layer provides the physical or virtual resources for slice deployment.

**Figure 3.11:** High level function of roles

3GPP in TR 28.801 [20] defines the following phases of a network slice lifecycle:

1. preparation phase;

2. instantiation, configuration and activation phase;

3. run-time phase;

4. decommissioning phase.

Moreover, it introduces three management functions to manage the NSIs to support communication services: Communication Service Management Function (CSMF), responsible for translating the communication service related requirements to network slice related requirements; the Network Slice Management Function (NSMF), responsible for management and orchestration of NSI; and the Network Slice Subnet Management Function (NSSMF), responsible for management and orchestration of a network slice subnet instance (NSSI). Finally, it defines the different roles of the actors (e.g., costumers, providers, operators, etc.) involved in slice provisioning as depicted in Fig. 3.11. However, TR 28.801 does not specify how to implement such functions and their relationship with respect to the ETSI NFV architectural framework.



**Figure 3.12:** Interaction between slice management functions (3GPP) and ETSI NFV Architecture

ETSI NFV EVE012 [21] establishes the correspondence between a network slice (3GPP) and a network service (ETSI NFV). There, ETSI describes that an NFV Network Service (NFV-NS) can be regarded as a resource-centric view of a network slice, for the cases where a NSI would contain at least one virtualized network function. Moreover, ETSI NFV EVE012 proposes that 3GPP slice management functions interact with ETSI NFV Architecture through the Os-Ma-Nfvo reference point as depicted in Fig. 3.12.

Within the research community, the H2020 Project 5G-Transformer (5GT) project [22] envisions three functional layers for providing verticals with slices: a

Vertical Slicer (5GT-VS) as the logical entry point for verticals to support the creation of their respective transport slices in a short time-scale (in the order of minutes), a Service Orchestrator (5GT-SO) to orchestrate the federation of transport networking and computing resources from multiple domains and manage their allocation to slices, and a Mobile Transport and Computing Platform (5GT-MTP), that provides and manages the virtual and physical IT and network resources on which slices are deployed. Such architecture implementation is under development. Chapter 8 focuses on slicing and in particular on the 5GT activities.

# 4 OpenAirInterface and the ARNO-5G Testbed

## 4.1 Introduction

In this section an overview of the software platform and hardware utilised to test and study the performance in different architecture deployment and uses cases of a Next Generation mobile network are presented.

## 4.2 OpenAirInterface Overview

In 1998 EURECOM [23] launched an experimental research activity for the development of next-generation wireless communications systems. This marked the birth of the convergence of wireless systems with the Internet. EURECOM quickly identified the need of a full-scale environment to test ideas for improvement of the components of a complete system. In this context the idea of federating efforts around the many required software components appeared necessary. The mission of the OpenAirInterface Software Alliance (OSA) [24] is therefore provide software and tools for 5G Wireless Research and Product Development.

The current generation of hardware/software for radio access network (RAN) consist of large numbers of proprietary elements that stifle innovation and increase the cost for the operators to deploy new services/application in an ever-changing fast paced cellular networks. Open source software running on general purpose processors (such as x86, ARM) can greatly simplify network access, reduce cost, increase flexibility, improve innovation speed and accelerate time-to-market for introduction of new services. There is already a movement going on within the industry on the development of Software Defined Networking (SDN) concepts to open the proprietary interfaces to control the RAN hardware/software. At the same time, open-source has made a very significant impact in the extremities of

**Figure 4.1:** OAI Architecture

current networks, namely in the terminals due to the Android ecosystem and in cloud infrastructure due, in part, to the OpenStack ecosystem.

An open source implementation of fully real-time stack (eNB, UE and core network) on general purpose processors when combined with SDN, Network Function Virtualization (NFV) and OpenStack can bring significant efficiency in RAN design from both innovation and cost perspective.

Established in 2014, the OSA is a French non-profit organization ("Fonds De Dotation"), funded by corporate sponsors.

OSA currently provides a standard-compliant implementation of a subset of Release 10 LTE for UE, eNB, MME, HSS, S-GW and P-GW on standard Linux-based computing equipment (Intel x86 PC/ARM architectures). The software is freely distributed by the Alliance under the terms stipulated by the OSA license model. It can be used in conjunction with standard RF laboratory equipment available in many labs (i.e. National Instruments/Ettus USRP and PXIe platforms) in addition to custom RF hardware provided by EURECOM to implement these functions to a sufficient degree to allow for real-time interoperation with commercial devices.

Some industrial users have already been working on OpenAirInterface (OAI)-based systems integrated with commercially-deployable remote radio-head equipment and have provided demonstrations at major industrial tradeshows (e.g., Mobile World Congress Asia 2014, Mobile World Congress Barcelona in 2013 and IMIC 2013). The primary future objective is to provide an open-source reference implementation which follows the 3GPP standardization process starting from Rel-13 and the evolutionary path towards 5G and that is freely-available for experimentation on commodity laboratory equipment.

The 5G cellular network evolution is under lot of debate amongst industry and academia. However, OSA have established several work areas which are regularly evaluated by the board of members in which alliance invests to ensure that OSA software meets the requirements of these strategic areas.

OSA broadly focuses on the evolution of 3GPP Cellular stack (eNB + UE + Core Network) on general purpose processor architectures (Intel/ARM) with the goal of establishing generic interfaces with 3rd party RF platforms like EURECOM Express MIMO, National Instruments/Ettus Research USRP, Nuand BladeRF, SoDeRa Lime SDR platforms.

The alliance also ensures the capability of running on Commercial-Off-The-Shelf (COTS) hardware platforms, for example Intel x86 and ARM. The Fig. 4.1 shows the conceptual architecture of OAI and how it relates to several hardware RF platforms.

## 4.3 ARNO-5G Testbed

A LTE/5G network exploiting the OSA software and tools has been set in the ARNO (Advanced Research on NetwOrking) testbed in the Scuola Superiore Sant'Anna - TeCIP Institute Laboratory [25].

In the 5G segment of the Advanced Research on Networking (ARNO-5G) testbed [26], several LTE/5G network configurations have been deployed exploiting the available general purpose processors. A block diagram of all the ARNO-5G testbed is shown in Fig. 4.2.

Each segment of the ARNO-5G testbed depicted in Fig. 4.2 has been deployed accordingly to the considered PC specifications. Indeed for the core network deployment (Core Segment in Fig. 4.2) low-performance PCs are sufficient contrarily to what is required for the DU/eNB/UE Segment. Table 4.1 recaps the major hardware attributes of the machines belong to the ARNO-5G testbed.

**Figure 4.2:** ARNO-5G testbed block diagram

It is noted that the CU segment is used only in a C-RAN deployment. Accordingly to the considered RAN deployment, in the DU/eNB/UE Segment, eNBs or DUs are deployed.

Regarding the RF front-end of the DU/eNB/UE segment, different type of devices have been considered to test several 3rd party RF platforms. The considered RF devices (e.g., Ettus USRPs and LimeSDR) performs Digital to Analog and Analog to Digital Conversion (DAC/ADC), Digital Up and Down Conversion (DUC/DAC), low pass filtering, and amplification. Table 4.2 summarizes the main features of the RF devices belong to the ARNO-5G testbed.

At the end of the LTE/5G chain, two Huawei E3372 LTE dongles and two Huawei P8 Lite smartphones are utilized as User Equipments (UEs). The Huawei E3372 dongles support LTE category 4 and frequency-division duplexing (FDD) communication in the following bands: 900 MHz, 1800 MHz, 2100 MHz and 2600 MHz. They support a maximum rate of 150Mb/s in downlink and 50Mb/s uplink with a signal bandwidth of 20MHz. The dongles are connected to the RF devices (Ettus USRPs and LimeSDR) through SMA cables with 40 dB of attenuation. The Huawei P8 lite smartphones support the LTE bands, 3G bands (e.g., 1900/2100/850/900 MHz), 2G Bands (e.g., 1800/1900/850/900 MHz) and the GPRS and EDGE bands. The Huawei P8 Lite smartphones can connect to the LTE/5G network through the air interface utilising preconfigured SIM cards registered to

the HSS deployed in the core network.

| PC ID | PC Type | Processor Type | OS |
|---|---|---|---|
| PC 1, PC 2, PC 3, PC 4 | Up-board I Generation | Intel Atom x5-Z8350 Quad Core @ 1.9 GHz | Ubuntu 14.04 (4.7 kernel) |
| PC 5, PC 6 | Dell T410 PowerEdge | Intel Xeon E5620 Quad Core @ 2.0 GHz | Ubuntu 14.04 (3.19 low-latency) |
| PC 7 | Mini-ITX | Intel I7 7700 Quad Core @ 4.0 GHz | Ubuntu 14.04 (3.19 low-latency) |
| PC 8 | Desktop Computer | Intel i7 4790 @ 3.60 GHz | Ubuntu 14.04 (3.19 low-latency) |
| PC 9, PC 10, PC 11 | Intel NUC 7 | Intel i7 7567U Quad Core @ 3.5 GHz | Ubuntu 14.04 (3.19 low-latency) |

**Table 4.1:** Major hardware attributes of the ARNO-5G testbed machines

All the PCs composing the ARNO-5G testbed, have a management plane interface and also a data plane interface. The management and data planes belong to two different networks. The former one (i.e., 10.30.x.x) belonging to the lab backbone and handled by a Cisco Catalyst 2960G switch (indicated as SWITCH 2 in Fig. 4.2), ensures the continuous reachability of the machines. The latter one is handled by a second Cisco Catalyst 2960G switch (indicated as SWITCH 1 in Fig. 4.2) and it is used exclusively as the ARNO-5G testbed data plane for the communication between the LTE network entities (i.e., OAI control and user plane). All the PCs are connected to the management plane through 100 Mb/s Ethernet links while they are all connected to the SWITCH 1 by a 1 Gigabit Ethernet link.

The SWITCH 1 is configured to have different subnets for the backhaul link, the midhaul link and the fronthaul link. For example, the EPC and the CU interfaces belonging to the backhaul link are configured in the 10.10.20.x subnet while the midhaul link and the interfaces of the CU and DU are configured in the 10.10.30.x subnet.

ARNO-5G testbed is federated in the Fed4FIRE federation and it accepts users from only one trusted central authority (iMinds) identity provider. Therefore, external experimenters can configure experiments interconnecting resources from multiple testbeds at the same time, reserve and access them via iMinds tools such as jFed [27]. Once successfully logged in, experimenters can set up their experiments by choosing which types of resources and from which testbeds, configure those resources (operative system, software to be installed, network configurations, measurement options, etc.), launch the experiments and access the

| ETTUS B210 | ETTUS X310 | LimeSDR |
|---|---|---|
| 2 channels | up to 4 channels | |
| 70 MHz-6 GHz | 50 MHz-2.2 GHz | 100 kHz-3.8 GHz |
| Up to 56 MHz of bandwidth | Up to 160 MHz of bandwidth | 61.44 MHz of bandwidth |
| Full duplex 2x2 MIMO | 2 RF daughter board slots | 2x2 MIMO |
| SuperSpeed USB 3.0 | Multiple high speed IF | SuperSpeed USB 3.0 |
| GNU Radio | GNU Radio | LimeSDR-USB |
| Xilinx Spartan6 | Xilinx Kintex7 | Altera Cyclone IV |
| AD9361 RFIC transceiver | 14 bit ADC 16 bit DAC | Lime Microsystems LMS7002M transceiver |

**Table 4.2:** ARNO-5G RF Devices specifications

resources. Through jFed tool, an experimenter can select the ARNO-5G testbed, namely "Sant'Anna Pisa testbed" and provide their slice name. This process creates a Docker container in ARNO-5G testbed that acts as a gateway towards the ARNO-5G testbed OAI components. Then experimenters can enter each OAI component of ARNO-5G testbed through ssh based on the specific container. More details about how to reserve the components of ARNO-5G testbed can be found in [26].

In all the experiments and studies explained in the following of this thesis the ARNO-5G testbed was used, considering in some case all the elements belonging to it and in some other case only a subset of them to implement different scenarios and different LTE/5G network configuration as described in the appropriate chapters.

# 5 NG RAN Virtualisation: Deployment and Performance Analysis

## 5.1 Introduction

In this chapter is presented an evaluation of how different virtualization technologies decrease the midhaul latency budget when different 5G network components are virtualised considering different gNB functional splits. The midhaul latency budget and packet jitter (i.e., packet delay variation) budget are computed in all the possible combinations and compared with the scenario when all the network components belonging the 5G network are deployed in bare metal.

The considered performance evaluation parameters are the Allowable Latency Budget (ALB) and the Allowable Jitter Budget (AJB) of the midhaul segment. In the experimental performance evaluation presented in the following sections, the contribution of the link connecting the RU and DU, known as fronthaul I interface, is assumed to be negligible because the RU is not virtualized and the link between RU and DU is assumed to be short. The ALB and the AJB are defined as the maximum one-way latency and the maximum latency variation (i.e., delay jitter) supported by the midhaul segment without disconnection. A disconnection occurs when the latency and the jitter of the midhaul segment cause loss of synchronization between CU and DU.

To emulate latency and jitter in the midhaul, the Linux traffic control (tc) tool is utilized. The tc utility is based on a token bucket filter and it is capable of increasing the delay and jitter experienced on a link by a packet by storing it in the output interface for a specified amount of time before its transmission on the link. A delay d0 is applied to the Ethernet interface of the machine in which the DU is deployed and a delay d1 is applied to the Ethernet interface of the machine in which the CU is deployed. In this way a one-way latency is inserted in the midhaul link. For evaluating the ALB, d0 and d1 are increased with steps of 10 $\mu$s until DU, CU, and UE disconnect. For evaluating the AJB, instead, the jitter follows a

normal distribution and it is added to the latency values d0 and d1 with steps of 10 $\mu$s until DU, CU, and UE disconnect. Two different scenarios are considered in the AJB evaluation:

1. In the first case, the mean latency is set to 95% of the ALB and the supplementary random delay is increased to discern if jitter could be the origin of an ALB reduction.

2. In the second case, the mean latency is set to 42.5% of the ALB and the supplementary random delay is increased to discern if jitter could be the limitation for the midhaul segment.

In particular, Section 5.2 overview possible gNB functional splits and illustrates which requirements the Radio Access Network (RAN) shall satisfy to effectively transport 5G protocol data units (PDUs) as a function of the considered functional split.

Section 5.3 instead presents a preliminary experimental evaluation of the midhaul performance in terms of ALB and bandwidth consumption (MB/s) when a single gNB functional split is implemented.

Moreover in Section 5.4 and 5.5, an evaluation of how different virtualization technologies (i.e., VirtualBox, Kernel-based Virtual Machine, and Docker Container) decrease the midhaul latency budget considering CU virtualization and Option 7-1 functional split is performed.

Finally Section 5.6 evaluates many additional scenarios. In particular, several virtualization technologies are utilized to virtualize not only the CU but also the DU. Both split Option 8 and Option 7-1 are considered. The midhaul latency budget and packet jitter (i.e., packet delay variation) budget are computed in all the possible combinations and compared with the scenario when CU and DU are deployed in bare metal. A mathematical model, expressing the midhaul latency budget as a function of the considered channel bandwidth, functional split options, and virtualization technologies is provided and validated through experimental results. In an additional experimental analysis, how a vDU performance is impacted by virtualized elements (i.e., CU and DU) deployed in the same computational resource is studied. In this way, the need for anti-affinity constraint when a vDU is deployed is evaluated. Finally, the impact of deploying several vDUs/vCUs in the same host (i.e., the VRAN scalability) is experimentally evaluated.

# 5.2 Performance Evaluation #1: Protocol Stack Layer overhead

In VRAN, CUs, based on Virtual Machines (VM), can be activated in general purpose servers located anywhere in the RAN. However how to transport both the Control plane (C-plane) and the User Plane (U-Plane) data between DUs and CUs remains an open issue as well as which functional split optimizes performance [28, 29].

A capacity requirements of different gNB functional splits must therefore be evaluated to understand if will be possible to satisfy the rigid objectives imposed by the 5G.

This section overviews the possible gNB functional split to show to the reader where the LTE protocol stack layers are placed in the network entities composing the RAN (e.g. RU, DU, CU). Therefore experiments are performed considering a classical gNB without functional split (i.e., all the LTE protocol layers are in the same "host") to calculate the C-plane and U-plane overhead and the overhead consumption capacity of each layer. Finally, the percentage contribution in terms of C-Plane and U-Plane overhead of each layer in the considered functional splits are measured.

The envisioned gNB functional split in CU and DU is depicted in Fig. 5.1. Four of the five layers of the LTE-A Pro gNB architecture are considered as potential splitting points: Medium Access Control (MAC), Radio Link Control (RLC), Packet Data Convergence Protocol (PDCP), and Radio Resource Control (RRC).

The midhaul transports the PDU of the lowest layer of the CU stack to the DU and vice versa. Each gNB functional split comes with different requirements for the midhaul transport protocol based on the functions each layer performs. A short description of the different functions performed by the gNB layers and the related requirements follows and it based on what is reported in [30].

The Radio Resource Control (RRC) sublayer main services and functions include: broadcast of System Information (i.e., System Information Blocks - SIB); paging; establishment, maintenance and release of an RRC connection between the User Equipment (UE) and Evolved Universal Terrestrial Radio Access Network (EUTRAN); security functions including key management; establishment, configuration, maintenance and release of point to point Radio Bearers (RB); mobility functions; notification and counting for Multimedia Broadcast Multicast Ser-

**Figure 5.1:** An example of considered functional splits for 5G systems

vice (MBMS) services; establishment, configuration, maintenance and release of Radio Bearers for MBMS services; QoS management functions; UE measurement reporting and control of the reporting; NonAccess Stratum (NAS) direct message transfer to/from NAS from/to UE. The Packet Data Convergence Protocol (PDCP) sublayer main services and functions for the user plane include: header compression and decompression (i.e., Robust Header Compression - ROHC); transfer of user data; in-sequence delivery of upper layer PDUs at PDCP re-establishment procedure for Radio Link Control (RLC) Acknowledged Mode (AM); PDCP PDU routing for transmission and PDCP PDU reordering for reception for split bearers in Dual Connectivity (DC); duplicate detection of lower layer SDUs at PDCP; retransmission of PDCP SDUs at handover and, for split bearers in DC, of PDCP PDUs and PDCP data-recovery procedure, for RLC AM; ciphering and deciphering; timer-based SDU discard in uplink. The PDCP sublayer main services and functions for the control plane include: ciphering and integrity protection; transfer of control plane data.

The Radio Link Control (RLC) sublayer main services and functions include: transfer of upper layer PDUs; error correction through ARQ (only for AM data transfer); reordering of RLC data PDUs; duplicate detection; protocol error detection; RLC SDU discard; RLC re-establishment.

The Medium Access Control (MAC) sublayer main services and functions include: mapping between logical channels and transport channels; multiplexing/de-multiplexing of MAC SDUs belonging to one or different logical channels into/from

transport blocks (TB) delivered to/from the physical layer on transport channels; scheduling information reporting; error correction through Hybrid Automatic Repeat request (HARQ); priority handling between logical channels of one UE; priority handling between UEs by means of dynamic scheduling; MBMS service identification; transport format selection; padding: The MAC sublayer includes also a sidelink whose main functions are: radio resource selection; packet filtering for sidelink communication.

Different gNB functional split will therefore bear different requirements for the midhaul both in terms of capacity and latency. Some papers already evaluated such requirements when the functional split is implemented in different sublayers of the physical and upper layers. In [31–33] the bandwidth and latency requirements of six different gNB functional splits are presented for LTE. The six proposed functional splits for the uplink direction are: fully centralized gNB, PHY1, PHY2, MAC-PHY, MAC-MAC, PDCP-RLC.

The fully centralized mapping is based on transporting, through Common Packet Radio Interface (CPRI), radio link time-domain in-phase and quadrature (I/Q) samples to the CU location.

The PHY1 functional split is based on maintaining at the RRH the function of serial to parallel conversion of the data, the removal of cyclic prefix and of the application of the Fast Fourier Transform (FFT). Thus, frequency-domain I/Q samples (the channels are still multiplexed) of all the cell resource blocks (RBs) are transported together with the Physical Random Access Channel (PRACH) data.

The PHY2 functional split still transports the frequency-domain I/Q samples but only the ones utilized by the User Equipments (UEs), that is after the resource mapping. Thus, the capacity depends on the actual load of the wireless system.

In the MAC-PHY functional split all the functionalities related to signal generation/detection are performed at the RRH. Thus uncoded user payloads in the form of MAC PDUs including all higher layer overheads are transported between RRH and BBU.

The MAC-MAC functional split introduces Transmission Selection Unit (TSU) whose task is to select one of the pre-calculated scheduling assignments from the CU based on results from the UL PHY and forward it further to the DL PHY to generate the next subframe accordingly.

In the PDCP-RLC split the RLC, MAC and PHY processing is performed remotely at the DU. Thus, PDCP frames only need to be transported between DU and CU saving the data overhead added by lower layers (e.g., MAC and RLC head-

ers are used at the remote units).

In terms of capacity requirements in [32] a scenario with 20 MHz channel width, 1200 subcarriers, 2x2 MIMO, and 50% Resource Block Utilization is considered. In [32] it is reported that a MAC-PHY functional split would require almost 1% of the fully centralized solution which requires about 2.5 Gb/s. In terms of latency requirements in [33] it is reported that latency requirements at the PHY layer are in order of 1ms, at the MAC layer are of about 8ms (due mainly to the HARQ RTT), at the RLC are in the order of hundreds of ms and at the RRC are in the order of units or tens of seconds.

In this section, the C-plane and U-plane overhead at each layer and C-plane functional split overhead as a function of the different functional split options is considered. The C-plane overhead at each layer is defined as the overall amount of control data and overhead exchanged in bytes including SIBs. The U-plane overhead at each layer is defined as number of overhead bytes used to transport the considered application data. The C-plane Functional Split Overhead (CFSO), is defined as the sum of C-plane overhead of each layer residing at the CU based on the implemented functional split option, as depicted in Fig. 5.1. The CFSO percentage is defined as the ratio between CFSO and the offered load. Moreover, the overhead capacity consumption is defined as the ratio between the overhead and the experiment duration for C-plane overhead, U-plane overhead, and CFSO.

The experimental evaluation of the capacity requirements of the gNB functional splits depicted in Fig. 5.1 is performed by means of OAISIM (OpenAirInterface System Emulation) [34]. OAISIM is considered to emulate the RAN and the UE. OAISIM, particular feature of the OAI emulation platform presented in Section 4 allow to simulate the LTE RAN segment. It provides either simulation with full physical layer (PHY) and synthetic radio channels, or using PHY abstraction. Without losing in generality, a single gNB and a single UE is considered, and the UE is located at a fixed distance D (i.e., 370 m) from the gNB. Only uplink transmission (i.e., UE to gNB) is considered throughout the experiment duration. OAI provides different traffic models with Constant Bit Rate (CBR) traffic patterns: Small-packet CBR (SCBR), Medium-packet CBR (MCBR) and big-packet CBR (BCBR). The SCBR and MCBR traffic patterns are more suitable for machine-to-machine type communication, whereas BCBR is suitable for conventional human-type communication. In this study the SCBR traffic pattern is considered to better understand the impact of C-plane overhead on the considered functional split options. Furthermore, the user datagram (UDP) protocol is con-

| Parameter | Value |
|---|---|
| Simulation Duration | 100000 TTIs |
| Duplexing Mode | FDD |
| PHY Layer Abstraction | NO |
| # gNBs | 1 |
| # UEs | 1 |
| Mobility | STATIC |
| Payload Size | 200 bytes |
| IDT | 1 ms |
| Offered load | 1.6 Mb/s |
| Traffic Type | SCBR |
| TX mode | 1 (SISO) |
| Carrier Bandwidth | 5 MHz |
| Multipath channel simulation | AWGN |

**Table 5.1:** Performance Evaluation # 1: Simulation Parameters

sidered as an ISO/OSI transport layer protocol. The simulation duration is set 100 s, which is a 105 ms transmission time intervals (TTI), and the distribution of inter-departure time (IDT) is set to uniform with duration of 1 ms. Unless stated otherwise, all the other simulation parameters are as shown in Table 5.1.

The Wireshark tool is used to analyse the packets exchanged between the gNB and the UE. OAI provides Layer 2 and above LTE protocol stack control and data signals tapping with help of Wireshark or packet capture (PCAP) trace. The following logical channels are considered from the captured trace: Broadcast control channel (BCCH), Common control channel (CCCH), Dedicated control channel (DCCH), periodic broadcast of system information for non-access stratum and access stratum (e.g., SIBs).

Among these control signals, only CCCH and DCCH signals are encapsulated into the lower layer PDUs (PDCP, RLC, and MAC layers). BCCH and (observed) SIBs (i.e., SIB1, SIB2, and SIB3) are directly encapsulated into MAC layer PDU. All these signals are encapsulated at MAC with transport channel (shared channel-SCH) based on its directions (e.g., DL-SCH or UL-SCH), as reported in [35].

Table 5.2 shows C-plane and U-plane overhead and overhead capacity consumption at each layer. In addition, it also shows different logical and transport channels utilized at each layer. As expected, RRC contribute more number of bytes

| Layer | Overhead [bytes] | | Overhead Capacity Consumption [b/s] | | Logical/Transport Channels |
|---|---|---|---|---|---|
| | C-Plane | U-Plane | C-Plane | U-Plane | |
| RRC | 473518 | - | 37881.44 | - | BCCH, CCCH, DCCH |
| PDCP | 115 | 45005 | 9.20 | 3600.40 | DCCH |
| RLC | 182 | 18002 | 14.56 | 1440.16 | DCCH, SRB1 |
| MAC | 231358 | 4112279 | 18508.64 | 328982.45 | BCH, SCH |

**Table 5.2:** Performance Evaluation # 1: Layer Overhead Results



**Figure 5.2:** Performance Evaluation # 1: gNB-UE messages exchange

to the C-plane overhead, because of the periodic transmission of SIB information as confirmed by the packet capture reported in Fig. 5.2 red rectangle. On the other end, PDCP and RLC have less contribution to C-plane overhead because these layers are involved only in transportation of DCCH for non-access stratum (NAS)-evolved packet core (EPC) connection. Moreover, RLC has other additional C-plane overhead to specify signalling radio bearers (SRB) for RRC and NAS signalling messages. The C-plane MAC overhead includes MAC PDU header, buffer status report (BSR), and power headroom report (PHR) of MAC control elements.

In general, the U-plane overhead depends on the total transport protocol headers (i.e., GPRS Tunnelling Protocol for U-plane (GTP-U), UDP, IP and Ethernet) with/without Internet Protocol Security (IPsec). As shown in Fig. 5.2 blue rectangle, the U-plane protocol headers with IPsec have a dimension of 144 bytes. The sum of all the considered U-plane overhead capacity consumption is about 21% of the offered load but this depends on the modulation and coding scheme (MCS) and the number of resource blocks assigned to the UE, which is a transport block size (TBS).

Fig. 5.3 shows the impact of functional split options as a function of CFSO.

**Figure 5.3:** Performance Evaluation # 1: CFSO measurements

The CFSO of MAC split has around 3.6% of offered load, and all other there splits are around 2.4% of offered load. However, this CFSO can change depending on application payload size if random payload size is assumed [36]. Moreover Fig. 5.3 results show that PDCP/RRC functional split requires low C-Plane overhead while bringing advantages in the terms of load balancing, mobility management and energy efficiency due to the possibility of virtualizing PDCP/RRC sublayer functions because of their loose latency constraints.

Based on the obtained results the capacity required for transporting both user and control data are limited to slightly less than one Megabit per second. Thus the limiting requirement for the midhaul is the latency which can vary from units of milliseconds, if a physical layer functional split is considered, to tens of seconds if an RRC layer functional split is implemented.

## 5.3 Performance Evaluation # 2: Latency limits of gNB functional split

One of the main constraints for functions placement of a gNB, is the latency experienced by the communication between the Virtual Machines (VMs) hosting the functions.

This section evaluates experimentally the performance, in terms of latency and bandwidth consumption (Mb/s) of the midhaul link, when a single gNB is splitted in a DU and a CU. The DU and CU are deployed in bare metal to derive a latency

**Figure 5.4:** Performance Evaluation #2: ARNO-5G Configuration

limits that must be respected in virtual environments.

To perform the experimental evaluation the 5G-ARNO testbed configuration depicted in Fig. 5.4 has been setup.

The EPC and the functional elements belonging to it (i.e., the Serving Gateway (S-GW), the PDN Gateway (PDN-GW), the Mobile Management Entity (MME) and the Home Subscriber Server (HSS)) are deployed on PC 1. The CU is deployed on PC 5 and it is connected by a 1 Gigabit Ethernet link to the EPC. The DU is deployed on PC 7 and it is connected to the CU by a 1 Gigabit Ethernet link as well. The DU is also connected through a USB 3.0 link to an Ettus B210 USRP. The UE is deployed by means of a Huawei E3372 LTE dongle. Please refer to Table 4.1 and 4.2 for the PCs and RF devices specifications considered in this section.

In such first simple scenario the functional split IF4p5 (also known as Option 7-1 in the 3GPP terminology) and signal bandwidths equal to 5 MHz and 10 MHz (corresponding to 25 and 50 Physical Resource Blocks - PRBs) are considered. In the UL direction of Option 7-1 functional split as defined by 3GPP, the UL, FFT, CP removal and possibly PRACH filtering functions reside in the DU and the rest of PHY functions reside in the CU. In the DL direction instead, iFFT and CP addition functions reside in the DU, the rest of PHY functions reside in the CU. In other word the Option 7-1 functional split is made before/after the resource mapping/demapping respectively. Without losing in generality, a single EPC, a single

**Figure 5.5:** Performance Evaluation #2: Midhaul bandwidth (5 MHz)

CU, a single DU and a single UE are considered. The UE is static and connected to the DU through coaxial cables. The experimental evaluation of the midhaul bandwidth occupation, is carried out in three different scenarios. In the first scenario no UEs are attached. In the second scenario a UE is attached but only control plane data are sent. In the third scenario, when the UE is attached, the offered traffic in the uplink is gradually increased. In particular, from the UE, ping tests with a fixed packet size, equal to 512 bits, and an incremental inter-departure time are performed. With the -i ping option, the waiting time between sending each packet is set. In our experiment the waiting time starting from 0.1 second is divided by $10^{-1}$ in each steps until the value $10^{-5}$ is reached.

In this way the offered traffic, calculated by means of Eq. 5.1 , varies from 5.12 kb/s to 5.12 Mb/s respectively.

$$Offered\ Traffic\ [b/s] = \frac{Ping\ Packet\ Size}{Ping\ Packet\ Waiting\ Time} \tag{5.1}$$

Fig. 5.5 and Fig. 5.6 show the midhaul bandwidth when a signal bandwidth equal to 5 MHz and 10 MHz is considered respectively. As expected, the midhaul bandwidth occupation is fixed and not sensible at the UE traffic. It only depends on the signal bandwidth and on the applied functional split. As depicted by Fig. 5.7 when the delay overcomes a certain delay threshold DU and CU are not capable of communicating. The threshold is about 200 $\mu$s for the 10 MHz bandwidth and about 250 $\mu$s for the 5 MHz bandwidth. The emulated delay causes loss of synchronisation between the DU and the CU. The module that performs the FFT/IFFT, implemented by the USRP Ettus B210 at the DU, receives samples from the CU

**Figure 5.6:** Performance Evaluation #2: Midhaul bandwidth (10 MHz)



**Figure 5.7:** Performance Evaluation #2: Midhaul latency limit

through the midhaul, not in the expected order. For these reasons a mismatch occurs and the connectivity between the CU Fig. 5.7 show the status of the midhaul connection between DU and CU as a function of the link latency and of the signal bandwidth.

From the obtained results is possible to conclude that the capacity requirement is independent of the traffic generated by the UE because the midhaul is carrying cell-level information. Moreover, the maximum one-way latency that can be tolerated along the midhaul is about 250 $\mu$s as specified by 3GPP. and the DU is lost. It must be noted that the delay midhaul threshold does not change significantly with the signal bandwidth.

## 5.4 Performance evaluation #3: Increasing the DUs

This section evaluates experimentally the ALB and AJB that Option 7-1 functional split can support in the midhaul using different radio channel bandwidths and when different number of DUs and different number of User Equipments (UEs) are considered in a physical environment.

The analysis presented in this section is performed by means of the ARNO-5G testbed [26] configuration shown in Fig. 5.8.

The EPC and the functional elements belonging to it (i.e., the Serving Gateway (S-GW), the Public Data Network Gateway (PDN-GW), the Mobile Management Entity (MME) and the Home Subscriber Server (HSS)) are deployed in PC 1 (please refer to Table 2.1 for the PC1 specifications).

The SWITCH 1 in Fig. 5.8 acts as Radio Aggregation Unit (RAU). The RAU becomes a necessary network element because of the point-to-multipoint architecture between the CU and the DUs. The RAU forwards the communication from the CU to several DUs.

The CU is deployed in PC 5 (please refer to Table 2.1 for the PC 5 specifications) and it is connected by a 1 Gigabit Ethernet link to the EPC and by a 1 Gigabit Ethernet link to the DU as well.

A first DU (DU1) is deployed in PC 7 (please refer to Table 2.1 for the PC 7 specifications). This machine is connected to the CU by a 1 Gigabit Ethernet link. It is also connected through USB 3.0 link to an Ettus B210 for implementing the Radio Frequency (RF) front-end. A second DU (DU2) is deployed in PC 8 (please refer to Table 2.1 for the PC 8 specifications) and is connected as the DU1 to the RAU through a 1 Gigabit Ethernet link. Also the DU2 is connected through USB 3.0 link to a second Ettus B210 for implementing the RF front-end.

The UEs (i.e., UE1 and UE2) consists of Huawei E3372 LTE dongles connected to PC 9 and PC 11 respectively. Please refer to Table 4.1 and 4.2 for the PCs and RF devices specifications considered in this section.

The ALB and AJB evaluations are performed for signal bandwidths equal to 5 MHz and 10 MHz, corresponding to 25 and 50 Physical Resource Blocks (PRBs).

TCP traffic is generated by using iperf tool to check the UE connectivity stability. The EPC acts as iperf server while the UE acts as iperf client.

Two different scenarios are considered as described as follows. In Scenario 1, a single DU (i.e., DU1) is connected to a single CU through the RAU. In this scenario, having only one DU, we bind a single interface with a single UDP port number to

**Figure 5.8:** Performance Evaluation #23: ARNO-5G Configuration

the DU 1. All the RAN and EPC functional elements are run on physical machines. In Scenario 2 the two DUs (i.e., DU1 and DU2) are connected with a single CU. In order to be able to create connectivity for both the DUs we bind a single interface at the CU by using two different UDP port numbers one for each DU to serve them at the same time.

All the RAN and EPC functional elements run in physical machines. The two DUs are running in two different physical machines, while two OAI CU instances, running in the same physical machine, are connected to the corresponding DUs.

In both the aforementioned scenarios the UEs are static and connected to the DU through coaxial cables with 40 dB attenuation. The other experimental parameters are shown in Table 5.3.

Fig. 5.9 shows the ALB for the considered Scenario 1 and Scenario 2 with different signal bandwidth values (i.e., 5 MHz and 10 MHz). Results show that in the both considered scenarios the ALB is always below the 250 $\mu$s one-way latency constraint specified by 3GPP. Moreover, it can be noticed that the ALB decreases if the signal bandwidth and if the number of DUs connected to the same CU increases. The dependence on the signal bandwidth is due to the heavier processing required by the higher number of utilized PRBs. The dependence on the number of CU is

| Parameter | Value |
|---|---|
| Experiment Duration | 100000 TTIs |
| Duplexing Mode | FDD |
| Frame Duration | 10ms |
| PHY Layer Abstraction | NO |
| # DUs | 1 |
| # UEs | 1 |
| Mobility | STATIC |
| TX mode | 1 (SISO) |
| Carrier Bandwidth | 5 MHz, 10 MHz |

**Table 5.3:** Performance Evaluation #3: Experimental parameters

similarly due to the higher number of processes running in the same machine.

Latency requirements for different functional splits to serve high capacity New RAN architecture have been specified in the 3GPP. However, it is not clear how different functional splits can be affected by jitter. Thus, the second set of experiments aims at investigating whether the jitter impacts the ALB found in the first set of experiments.

In the considered experiments, we vary the jitter while keeping the midhaul latency fixed and within the above ALB. Fig. 5.10 shows the obtained jitter results in Scenario 1 and Scenario 2. In particular, in Scenario 1 the latency is set to 220 $\mu$s for the 5 MHz and is set to 160 $\mu$s when a 10 MHz signal bandwidth is considered. The obtained AJB in this case is equal to 35 $\mu$s and 30 $\mu$s for the 5 MHz and 10 MHz signal bandwidth, respectively. For Scenario 2, the experiments are carried out by setting a fixed latency on the midhaul link equal to 200 $\mu$s and 120 $\mu$s for 5 MHz and 10 MHz, respectively. The AJB is equal to 30 $\mu$s for the 5 MHz and 25 $\mu$s for the 10 MHz as depicted in Fig. 5.11.

By comparing the results reported in Fig. 5.9 and in Fig. 5.10 it can be deducted that jitter negligibly impacts the ALB. To observe the impact of the sole jitter on the midhaul link, that is to find the AJB, the latency value is set far from the ALB depicted in Fig. 5.9. The obtained results are shown in Fig. 5.11 for Scenario 1 and Scenario 2. In both Scenario 1 and Scenario 2, the latency is set to 100 $\mu$s and 50 $\mu$s for signal bandwidths 5 MHz and 10 MHz, respectively. In Scenario 1, the obtained AJBs are 30 $\mu$s and 25 $\mu$s for 5 MHz and 10 MHz signal bandwidth, respectively. Whereas, in Scenario 2, the obtained AJBs are 35 $\mu$s and 40 $\mu$s for 5 MHz and 10 MHz signal bandwidth, respectively.

**Figure 5.9:** Performance Evaluation #3: ALB



**Figure 5.10:** Performance Evaluation #3: AJB with a latency value close to the
ALB

**Figure 5.11:** Performance Evaluation #3: AJB with a latency far to the ALB

From the presented results, we can observe that when the jitter overcomes a certain threshold DU and CU are not capable of communicating. Indeed, the jitter cannot be higher than 40 $\mu$s because, if the jitter is large, the are periods in which not enough samples (i.e., modulation symbols) can be delivered to the PHY layer.

So, in conclusion the presented results shown that by increasing the instances of CU running in the same machine the ALB decreases of some tens of microseconds due to the higher number of computations required in the same machine. Similarly, but in the order of more than fifty microseconds, it happens if the signal bandwidth increases. Regarding the jitter instead, the evaluation results shown that jitter negligibly impact the ALB. However, the AJB is in the order of tens of microseconds in all the considered scenarios.

## 5.5 Performance evaluation #4: Moving to a virtual environment

Another foreseen feature of the ARNO-5G testbed is the possibility to virtualize different Radio Access Network (RAN) and Evolved Packet Core (EPC) function. Such feature is therefore exploited in this section to calculate the virtualized RAN and EPC limits and compare them with the deployment in physical machines.

**Figure 5.12:** Performance Evaluation #4: ARNO-5G Configuration

The midhaul latency and jitter evaluations are performed for different virtualisation softwares and for signal bandwidths equal to 5 MHz and 10 MHz, corresponding to 25 and 50 Physical Resource Blocks (PRBs). In this section, the considered experimental evaluation scenario is shown in Fig. 5.12.

Regardless of the considered hypervisor, in such scenario, the EPC and CU are virtualised while the DU is deployed in a physical machine. In particular, the EPC is deployed in PC 1. The Mobile Management Entity (MME) is deployed in a VM called MME-VM and the Home Subscriber Server (HSS) is deployed in a second VM called HSS-VM. The Serving Gateway (S-GW) and the PDN Gateway (PDN-GW) are deployed in a third VM called SPGW-VM. The CU is hosted in PC 5 and is deployed in another VM called CU-VM. The DU runs always in PC 7 in a physical machine and the UEs, connected to the RAN through SMA cables with 40 dB of attenuation, are deployed by means of a Huawei E3372 dongle connected to PC 8. Please refer to Table 4.1 and 4.2 for the PCs and RF devices specifications considered in this section.

Different virtualisation methods are considered: VirtualBox, Kernel-based Virtual Machine (KVM) and Docker. Using VirtualBox, MME-VM, HSS-VM and SPGW-VM host Ubuntu 16.04 with 4.8 generic kernel featuring a one core virtual

CPU and 1 GB of RAM. Instead CU-VM hosts Ubuntu 14.04 with 3.19 low-latency kernel featuring a 8 core virtual CPU and 16 GB of RAM.

The utilized virtual Network Interface Controller (NIC) modes considered are bridged and host-only networking. Here, the bridge networking mode allows a VM to intercept data from/to the physical network effectively by creating a new network interface in software. Therefore we bridge a virtual ethernet interface in bridge networking mode in both MME-VM and SPGW-VM to the physical ethernet interface in PC 1, referred as br0 in Fig. 5.12 to implement midhaul communication. The host-only networking mode, is instead a networking mode that can be thought of as a hybrid bridged networking: the virtual machines can communicate to each other and the host as if they were connected through a physical Ethernet switch but they cannot communicate to external hosts since there is not a networking interface. Therefore, such networking mode is used for the internal communications between the entities composing the EPC: a virtual interface, different from the above mentioned, is set on MME-VM, HSSVM and SPGW-VM allowing the communications between the MME-VM and the HSS-VM (S6a interface) and between the MME-VM and SPGW-VM (S11 interface) through the hostonly adapter vboxnet0.

Whereas in CU-VM, two virtual interfaces are bridged in bridge networking mode with corresponding physical ethernet interfaces in PC 2 (referred as br0 and br1 in Fig. 5.12). This because the first CU-VM virtual interface has to connect to the MME-VM for the LTE control plane communications (S1-C interface) and with the SPGW-VM for the LTE data plane communications (S1-U interface). Instead, the second CU-VM virtual interface is used for the midhaul communication with the DU.

Summarizing, the NICs used for the virtualised EPC are:

- Bridge adapter enp0s3 to physical interface eno2 with subnet 10.30.x.x (used for the management plane);

- Host-only adapter enp0s8 to vboxnet0 interface, with subnet 192.168.x.x (used for internal EPC service configuration and relationship);

The NICs used for the virtualised CU are:

- Bridge adapter enp0s3 to physical interface eth5 with subnet 10.30.x.x (used to set S1-C and S1-U interface between CU and MME and between the CU and SPGW respectively);

- Bridge adapter enp0s9, to Physical interface eth1 with subnet 10.10.x.x (used to set the midhaul between DU and CU).

The second virtualisation software used for the deployment of the scenario shown in Fig. 5.12 is Kernel-based Virtual Machine (KVM). KVM, an open source software, is a full virtualisation solution for Linux on x86 hardware containing virtualisation extensions (Intel VT or AMD-V). It consists of a loadable kernel module that provides the core virtualisation infrastructure and a processor specific module. Using KVM, it is possible to run multiple virtual machines running unmodified Linux or Windows images. Each virtual machine has private virtualised hardware: a network card, disk, graphics adapter, etc. For our purpose we characterised the MME-VM, HSS-VM and SPGW-CU with Ubuntu 16.04 with 4.8 generic kernel featuring a one core virtual CPU and 1 GB of RAM. Instead CU-VM hosts Ubuntu 14.04 with 3.19 low-latency kernel featuring a 8 core virtual CPU and 16 GB of RAM. In PC 1 we bridge the management physical interface with a first interface of each VMs in pass-through source mode and the physical data plane interface with a second interface of each VMs in bridge source mode. In the pass-through source mode option, a virtual function of a Single-Root Input/Output Virtualisation (SRIOV) capable Network Interface Controller (NIC) is attached directly to a target VM without losing the migration capability. Therefore all the packets are sent directly to the network devices. In the bridge source mode option, packets whose destination is on the same host physical machine where they are originated from are directly delivered to the target device. Regarding the CU, because no internal host communication is needed, we set three different virtual interfaces in pass-through source mode each connected to three different physical interfaces. Thus, the NICs used for virtualised EPC are:

- pass-through adapter enp0s3 to physical interface eno2 with subnet 10.30.x.x (used for the management plane);

- bridge adapter enp0s8 to vboxnet0 interface with subnet 192.168.x.x (used for internal EPC service configuration and relationship);

The NICs used for the virtualised CU are:

- pass-through adapter enp0s3 to physical interface eth5 with subnet 10.30.x.x (used to set S1-C and S1-U interface between CU and MME and between the CU and SPGW respectively).

- pass-through adapter enp0s9 to physical interface eth1 with subnet 10.10.x.x (used to set the midhaul between DU and CU).

Finally, Docker has also been used for the deployment of the scenario shown in Fig. 5.12. Docker is an open platform for developers and it is a mechanism that helps in isolating the dependencies per each application by packing them into containers. Containers are more scalable to deploy than virtual machines. Virtual machines have a full OS with its memory management installed with the associated overhead of virtual device drivers. Containers are therefore smaller than Virtual Machines and enable faster start up with better performance, less isolation and greater compatibility which is possible due to sharing of the hosts kernel. Docker containers can share a single kernel and share application libraries. Containers present a lower system overhead than Virtual Machines and the performance of the application inside a container is almost the same as compared to the same application running on a Physical Machine but better as compared to Virtual Machine. There are different types of network modes available to connect Docker container with the host machine or to an external host. We connect our container through host network which uses the same protocol stack as the host device is using. A considerable advantage of using Docker containers is that we can bypass the overhead of bridge adapter used in previous approaches. Because OAI set many system variable values at run time, we run Docker container with privileged mode so that the container has got write permissions to set system variables. The following interfaces ((i.e., host's network interfaces) are set up in the utilized Docker containers:

- Physical interface eth5 with subnet 10.30.x.x (to set S1-C and S1-U interface between CU and MME and between the CU and SPGW respectively).

- Physical interface eth1 with subnet 10.10.x.x (used to set the midhaul between DU and CU)

Fig. 5.13 shows the midhaul ALB for the considered virtualisation methods with different signal bandwidth values (i.e., 5 MHz and 10 MHz). Using VirtualBox, the midhaul ALB is $40\mu$s for 5 MHz signal bandwidth. For 10 MHz signal bandwidth, CU and DU never communicate. This is due to the large number of samples generated at DU which cannot be handled with current configuration of the considered PC 2 in which CU-VM is deployed. If KVM is used, the midhaul ALB is 190 $\mu$s for 5 MHz bandwidth and 140 $\mu$s for 10 MHz bandwidth, respectively. By using

**Figure 5.13:** Performance Evaluation #4: Midhaul ALB

Docker, the midhaul ALB is 35 $\mu$s in the case of 5 MHz bandwidth and 165 $\mu$s for the MHz bandwidth. Thus, if VirtualBox is used, the midhaul ALB is very low when compared to when the CU is deployed in other virtualisation methods. By using KVM and Docker the ALB is close to the 3GPP constraint specified in TR 38.801. This is mainly due to how the packets are forwarded by the host network interface to the virtualized one and how they are managed. Regardless of the utilized virtualisation methods, the midhaul latency budget is also a function of the signal bandwidth. Such dependence is due to the heavier processing required by the higher number of PRBs.

Fig. 5.14 shows the obtained midhaul AJB results using the three virtualisation methods as above when the jitter is applied to a latency value close to the midhaul ALB.

With VirtualBox, the experiments are carried out by setting a fixed latency on the midhaul link equal to 20 $\mu$s for 5 MHz signal bandwidth. The obtained midhaul AJB is 25 $\mu$s and no communication was observed in case of 10 MHz signal bandwidth. In the KVM case, the experiments are carried out by setting a fixed latency on the midhaul link equal to 170 $\mu$s for a 5 MHz signal bandwidth and equal to 120 $\mu$s for a 10 MHz signal bandwidth. The midhaul AJB is equal to 20 $\mu$s for both signal bandwidths as shown in Fig. 5.14. With the Docker, the experiments are carried out by setting a fixed latency on the midhaul link equal to 220

**Figure 5.14:** Performance Evaluation #4: Midhaul AJB with latency close to the ALB

$\mu$s for 5 MHz signal bandwidth and equal to 150 $\mu$s in case of 10 MHz bandwidth. The obtained midhaul AJB is 25 $\mu$s in the first case and 35 $\mu$s in the second one. Thus, for all the considered virtualisation methods the latency budget is negligibly impacted by the jitter.

To observe the impact of jitter on the midhaul link, the latency value is set far from the budgets reported in Fig. 5.13. The obtained results are shown in Fig. 5.15 for all the considered virtualisation methods.

When the CU is virtualised with VirtualBox, the experiments are conducted by setting a fixed latency on the midhaul link equal to 20 $\mu$s for 5 MHz signal bandwidth. The obtained midhaul AJB is 25 $\mu$s and no communication was observed in case of 10 MHz signal bandwidth. Whereas, in KVM case, the latency value is fixed to 100 $\mu$s for 5 MHz and equal to 80 $\mu$s for 10 MHz signal bandwidth. The obtained midhaul AJB is 25 $\mu$s for the 5 MHz signal bandwidth and 30 $\mu$s for the 10 MHz case. With the Docker, the experiments are carried out by setting a fixed latency on the midhaul link equal to 150 $\mu$s for 5 MHz signal bandwidth and equal to 100 $\mu$s in case of 10 MHz bandwidth. The obtained midhaul AJB is 40 $\mu$s and 35 $\mu$s for 5 MHz and 10 MHz, respectively. Thus the maximum AJB is achieved by utilizing the Docker technology and it is about 40 $\mu$s. Fig. 5.15 shows the midhaul AJB with a fixed latency far from the ALB

**Figure 5.15:** Performance Evaluation #4: Midhaul AJB with a fixed latency far from the ALB

Results showed that lighter virtualisation methods (e.g., Docker) are impacting the midhaul latency budget for Option 7-1 (i.e., intra-PHY) split less than heavier virtualisation methods (e.g., VirtualBox). However, in all the cases, the midhaul latency budget reduction depends on the considered signal bandwidth. The higher the bandwidth the higher the computations required the higher the midhaul latency budget reduction. Furthermore, the performed experimental evaluation showed that a jitter of at most 40 $\mu$s can be tolerated.

## 5.6 Performance evaluation #5: Virtualisation Impact

The experimental analysis is carried out in the ARNO-5G [26] testbed configuration shown in Fig. 5.16.

In all the analyses presented in this section and described here below, the scenarios summarized in Tab. 5.4 are considered while the following key points are common to all the performed analyses:

- Three virtualization technologies are considered: Docker Container, Kernel-based Virtual Machine (KVM), and VirtualBox (VB);

- Two channel bandwidths i.e., 5 MHz and 10 MHz are considered.

**Figure 5.16:** Performance Evaluation #5: ARNO-5G Configuration

- Option 8 and Option 7-1 fucntional splits are considered. Both functional splits are Physical layer functional splits, as depicted in Fig. 5.17. and for both functional splits the midhaul latency requirement is about $250\mu$s one way, as specified in 3GPP TR 38.801.

In a first experimental analysis, ALB and AJB are measured by deploying the block diagram shown in Fig. 5.18 in the ARNO-5G testbed configuration shown in Fig. 5.16. The EPC is deployed in PC1. The CU is deployed in PC3 either in bare metal (i.e., CU) or virtualized (i.e., vCU), and one UE is considered. Similarly,

| Scenarios | Bare Metal | Virtualisation Technologies (Docker, KVM, VB) |
|---|---|---|
| Scenario 1 (S1) | DU and CU | x |
| Scenario 2 (S2) | DU | vCU |
| Scenario 3 (S3) | CU | vDU |
| Scenario 4 (S4) | x | vDU and vCU |

**Table 5.4:** Performance Evaluation #5: Experimental Scenarios

**Figure 5.17:** Performance Evaluation #5: Functional Split Options



**Figure 5.18:** Performance Evaluation #5: Scenario 1

the DU is deployed in bare metal or virtualized in PC5. The selected combination depends on which of the four aforementioned scenarios and summarized in Tab. 5.4, is considered. Based on the considered virtualization technology, the vCU and the vDU are installed in a Docker Container or virtual machine (VM).

A second experimental analysis aims at understanding if an anti-affinity constraint is necessary when multiple virtualized mobile functions with different functional split options are deployed in the same host. The anti-affinity constraint forces Virtualized Network Functions (VNFs) to be deployed in different computational resources. To perform such analysis, a more complex network is deployed by doubling the involved NG-RAN components as shown in Fig. 5.19.

An EPC, two CUs, two DUs, two RUs, and two UEs are deployed. The EPC is deployed in PC1. Either two bare metal processes of the CU or two vCUs are deployed in PC3. Similarly, either two bare metal processes of the DU or two vDUs are deployed in PC5. The vCUs and the vDUs are installed in a Docker Container or VM according to the considered virtualization technology. The bare

**Figure 5.19:** Performance Evaluation #5: Scenario 2



**Figure 5.20:** Performance Evaluation #5: Scenario 3

metal processes or the virtualized components are activated according to scenarios summarized in Tab. 5.4. In such analysis, three cases are examined:

(i) Only the Option 7-1 functional split is implemented.

(ii) The Option 7-1 functional split is implemented between a CU-DU pair and the Option 8 functional.

(iii) Only the Option 8 functional split is considered. split is implemented between the second CU-DU pair.

In a last analysis, the scalability of the system is verified by increasing the deployed NG-RAN components up to four as shown in Fig. 5.20. The experiment is conducted as in the first considered experiment.

Fig. 5.21 depicts the obtained ALB in the first experimental setup. The ALB obtained in S1 (i.e., the bare metal scenario) is considered as benchmark. As shown, the utilization of the Docker Container allows to reach the highest ALB in all the considered scenarios with virtualized elements because dockers are a lightweight virtualization technology.

Indeed, Docker Containers are a native application with respect to the host. Thus, they have a smaller footprint than the VMs implemented by means of KVM and VB. Furthermore, in KVM and VB, I/O virtualization is performed by means of a hardware emulation layer under the control of the hypervisor, introducing additional delay. In addition, ALB heavily depends on the channel bandwidth: wider channel bandwidths mean a larger number of Physical Resource Blocks (PRBs),

**Figure 5.21:** Performance Evaluation #5: Midhaul ALB

thus a high computing effort and a growing processing time are needed. In S2, if KVM is utilized, the ALB is zero when Option 8 is considered. Instead, with VB the UE is able to connect only when Option 7-1 and a 5 MHz channel bandwidth is used. In S3 and S4, ALB values are greater than zero with Docker Container only.

From the results reported in Fig. 5.21, it is possible to obtain an empirical formula that relates the ALB to the considered channel bandwidth, the functional split options, and the utilized virtualization technologies. Based on [37] and [38], the ALB can be expressed as:

$$ALB = T_{TH}^{3GPP} - T_{proc},$$
(5.2)

where $T_{TH}^{3GPP}$ is the midhaul latency threshold and the $T_{proc}$ is the sum of processing time at the DU and the CU. Based on the experimental results $T_{proc}$ can be linearly fitted as $T_{proc} = \alpha x + \beta$ where $x$ is the considered channel bandwidth and $\alpha$ and $\beta$ are coefficients depending on the virtualisation technology and split option. Tab. 5.5 shows the $\alpha$ and $\beta$ values estimated in the ALB experimental analysis performed by using the testbed shown in Fig. 5.18.

Fig. 5.22 (top) shows the ALB trend in S1. The results depicted in Fig. 5.22 (middle) and Fig. 5.22 (bottom) are obtained in S2 with Docker Container and KVM, respectively.

| Platform | Option 7-1 | | Option 8 | |
|---|---|---|---|---|
| | $\alpha$ | $\beta$ | $\alpha$ | $\beta$ |
| BM | -7.6 | 276 | -22.8 | 263 |
| Docker | -12.6 | 301 | -24.4 | 266 |
| KVM | -13 | 247 | x | x |

**Table 5.5:** Performance Evaluation #5: $\alpha$ and $\beta$ coefficients



**Figure 5.22:** Performance Evaluation #5: ALB trends as a function of channel bandwidth

**Figure 5.23:** Performance Evaluation #5: Midhaul AJB with a fixed latency close from the ALB

Note that only the Docker Container and KVM coefficients are calculated in S2 because the ALB values obtained when VB is used are too small to obtain a good fitting. As shown in Fig. 5.22, if split Option 7-1 is considered, Docker Container and KVM perform similarly as a function of the channel bandwidth (their $\alpha$ values are similar). As well, if split Option 8 is considered, BM and Docker Container perform similarly as a function of the channel bandwidth. These results confirm the capabilities of the Docker Container to achieve performance close to the ones of bare metal when the CU is virtualized.

Fig. 5.23 shows the AJB when the fixed mean latency is set equal to 95% of the ALB. The AJB values obtained in S1 are considered as benchmark. Fig. 5.24 shows the obtained results when the fixed mean latency is set equal to the 42.5% of the ALB.

The results show that in the former case even a small jitter can cause a disconnection between the CU and DU. In the latter case, higher AJB is allowed. In both considered cases the AJB is in the interval between $20\mu$s and $40\mu$s. Thus the midhaul is very sensitive to jitter. In all the considered scenarios for split Option 8 AJB is zero.

In S2 with VB, the UE is able to connect considering Option 7-1 and 5 MHz channel bandwidth only. In S3 and S4, AJB values are obtained with Docker Container

**Figure 5.24:** Performance Evaluation #5: Midhaul AJB with a fixed latency far from the ALB

and Option 7-1 split only.

Fig. 5.25 shows the results obtained in the second experimental setup (i.e., anti-affinity experimental analysis). The obtained results show that, in the case of 5 MHz channel bandwidth, if split Option 8 and Option 7-1 coexist in the same computational resource (labeled as Mixed), the ALB decreases with respect when only Option 7-1 is considered (note that in this case the ALB is the time of the first UE disconnection, that is the disconnection of the UE whose data plane utilizes split Option 8). In case of the 10 MHz channel bandwidth, in the considered setup, split Option 8 and split Option 7-1 cannot be deployed together (the achieved ALB is zero). In addition, it is not possible to deploy two split Option 8 with 10 MHz channel bandwidth in the same computational resource. Thus, anti-affinity constraint shall be imposed if VNFs featuring split Option 8 are deployed to avoid that split Option 8 ALB is heavily impaired.

For KVM and VB technologies UEs and DUs are not capable to communicate.

Fig. 5.26 shows the third experimental setup results (i.e., the scalability experimental analysis). Since the Docker Container resulted the best one among the analyzed virtualization technologies, the scalability experimental analysis is performed only with Docker Container. Results show that, the ALB decreases if the number of virtualized DU-CU components increases due to the greater traffic load

**Figure 5.25:** Performance Evaluation #5: Midhaul ALB in the anti-affinity con-
straint analysis

injected in the midhaul segment.

In this section an experimental analysis of the effect of virtualizing NG-RAN
components (e.g., CU and DU) on the maximum latency and jitter that the midhaul
can support has been presented. The first set of results showed that by using heav-
ier virtualization technologies and a higher number of physical resource blocks
(i.e., channel bandwidth), the midhaul maximum latency decreases due to a heav-
ier elaboration requested to the hardware. An empirical equation expressing the
midhaul maximum latency as a linear function of the number of physical resource
blocks (i.e., channel bandwidth), functional splits, and virtualization technologies
confirm the aforementioned trends. Moreover, even the midhaul jitter can be crit-
ical if it reaches values above $40\mu$s. A second set of results showed that if virtual
DUs and CUs featuring split Option 8 are deployed, the utilization of the anti-
affinity constraint is advisable to avoid large impairment in terms of maximum
supported latency. A third set of results showed that by increasing the number of
NG-RAN components in the same computational resource the maximum midhaul
latency heavily decreases.

**Figure 5.26:** Performance Evaluation #5: Midhaul ALB in the scalability analysis

## 5.7 Conclusion

In this chapter is presented an evaluation of how different virtualization technologies decrease the midhaul latency budget when different 5G network components are virtualised considering different gNB functional splits. In particular, Section 5.2 described the possible functional splits of the LTE-A Pro gNB between RRH and BBU. Then it analysed, through emulation, the capacity requirements for the C/VRAN midhaul. Based on the obtained results the capacity required for transporting both user and control data are limited to slightly less than one Megabit per second. Thus the limiting requirement for the midhaul is the latency which can vary from units of milliseconds, if a physical layer functional split is considered, to tens of seconds if an RRC layer functional split is implemented.

In Section 5.3 is presented an experimental evaluation of the capacity and the latency requirements of the midhaul network when the Option 7 (intra-PHY) functional split of the New Radio is implemented. As expected, the capacity requirement is independent of the traffic generated by the UE because the midhaul is carrying cell-level information. Moreover, the maximum one-way latency that can be tolerated along the midhaul is about $250\mu$s as specified by 3GPP.

Moreover, in Section 5.4 is presented the experimental evaluation of the impact of virtualizing gNB functions on the midhaul latency budget. It also showed the

maximum sustainable midhaul jitter. Results showed that by increasing the instances of CU running in the same machine the allowable midhaul latency budget decreases of some tens of microseconds due to the higher number of computations required in the same machine. Similarly, but in the order of more than fifty microseconds, it happens if the signal bandwidth increases. Moreover, if gNB functions are run in virtual machines the allowable latency budget further decreases, in the order of hundreds of microseconds, due to the higher number of computations required by the virtualization engine. Finally, the midhaul jitter evaluation showed that jitter negligibly impact the allowable latency budget. However, the allowable jitter budget is in the order of tens of microseconds in all the considered scenarios.

In Section 5.5 is instead presented the experimental evaluation of the impact of virtualizing gNB functions on the midhaul latency and jitter budget when different virtualisation methods are utilized. Results showed that lighter virtualisation methods (e.g., Docker) are impacting the midhaul latency budget for Option 7-1 (i.e., intra-PHY) split less than heavier virtualisation methods (e.g., VirtualBox). However, in all the cases, the midhaul latency budget reduction depends on the considered signal bandwidth. The higher the bandwidth the higher the computations required the higher the midhaul latency budget reduction. Furthermore, the performed experimental evaluation showed that a jitter of at most 40 $\mu$s can be tolerated.

Finally, in Section 5.6 an experimental analysis of the effect of virtualizing NG-RAN components (e.g., CU and DU) on the maximum latency and jitter that the midhaul can support has been performed. The first set of results showed that by using heavier virtualization technologies and a higher number of physical resource blocks (i.e., channel bandwidth), the midhaul maximum latency decreases due to a heavier elaboration requested to the hardware. An empirical equation expressing the midhaul maximum latency as a linear function of the number of physical resource blocks (i.e., channel bandwidth), functional splits, and virtualization technologies confirm the aforementioned trends. Moreover, even the midhaul jitter can be critical if it reaches values above $40\mu$s. A second set of results showed that if virtual DUs and CUs featuring split Option 8 are deployed, the utilization of the anti-affinity constraint is advisable to avoid large impairment in terms of maximum supported latency. A third set of results showed that by increasing the number of NG-RAN components in the same computational resource the maximum midhaul latency heavily decreases.

# 6 Multi-access Edge Computing (MEC) in 5G Networks

## 6.1 Introduction

Edge computing as an evolution of cloud computing brings application hosting from centralized data centres down to the network edge, closer to consumers and the data generated by applications. Edge computing is acknowledged as one of the key pillars for meeting the demanding Key Performance Indicators (KPIs) of 5G, especially as far as low latency and bandwidth efficiency are concerned. However, not only is edge computing in telecommunications networks a technical enabler for the demanding KPIs, it also plays an essential role in the transformation of the telecommunications business, where telecommunications networks are turning into versatile service platforms for industry and other specific customer segments. This transformation is supported by edge computing, as it opens the network edge for applications and services, including those from third parties.

ETSI ISG MEC (Industry Specification Group for Multi-access Edge Computing) has already published a set of specifications (Phase 1) focusing on management and orchestration (MANO) of MEC applications [39, 40], application enablement API [41], service Application Programming Interfaces (APIs) [42–45] and the User Equipment (UE) application API [46]. The MANO and application enablement functions contribute to enabling service environments in edge data centres, while the service APIs enable the exposure of underlying network information and capabilities to applications.

One of the key value-adding features of the MEC specification is this ability for applications to gain contextual information and real-time awareness of their local environment through these standardized APIs. This local services environment is a flexible and extendable framework, as new services can be introduced by following the API guidelines in [47], when creating new service APIs. And finally, the UE application API lets the client application in the UE interact with the MEC

system for application lifecycle management.

5G networks based on the 3GPP 5G specifications [48] are a key future target environment for MEC deployments. The 5G system specification and its Service Based Architecture (SBA) leverage the service based interactions between different network functions, aligning system operations with the network virtualization and Software Defined Networking (SDN) paradigms. These very same characteristics are shared by MEC specifications. In addition, 3GPP 5G system specifications define the enablers for edge computing, allowing a MEC system and a 5G system to collaboratively interact in traffic routing and policy control related operations. MEC features together with these complementary technical enablers of the 5G system allow integration of these systems to create of a powerful environment for edge computing. In Section 6.2, is presented a possible MEC integration in the 5G system solution and is described some possible MEC deployment solutions. In Section 6.3 the MEC contribution in the 5G Automotive Scenario is then introduced. Finally, in Section 6.4 and Section 6.5 two deployed MEC-based applications are presented.

## 6.2  MEC Deployment in next Generation 5G networks

MEC as it is deployed currently in the $4^{th}$ generation LTE networks, is connected to the user plane via one of the options described in the ETSI White Paper "MEC deployments in 4G and evolution towards 5G" [48]. With LTE networks already having been deployed for a number of years, it was necessary to design the MEC solution as an add-on to a 4G network in order to offer services in the edge. Consequently, the MEC system as defined in [49] and in the related interface specifications, is to a large extent self-contained, covering everything from management and orchestration down to interactions with the data plane for steering specific traffic flows.

With 5G, the starting point is different, as edge computing is identified as one of the key technologies required to support low latency together with mission critical and future IoT services. This was considered in the initial requirements. The system was designed from the beginning to provide efficient and flexible support for edge computing to enable superior performance and quality of experience.

The design approach taken by 3GPP allows the mapping of MEC onto Application Functions (AF) that can use the services and information offered by other 3GPP network functions based on the configured policies. In addition, a number

of enabling functionalities were defined to provide flexible support for different deployments of MEC and to support MEC in case of user mobility events. The new 5G architecture is described and explained in more detail in the following of this section.

## 6.2.1 5G & MEC System architectures

The 5G system architecture specified by 3GPP and described in [47] has been designed to cater for a wide set of use cases ranging from a massive amount of simple IoT devices to the other extreme of high bit rate, high reliability mission critical services. Supporting all the use cases with the same and common architecture has required significant changes in design philosophies both for the RAN and the core network. One significant architectural change was made to the communications between the core network functions that until now have relied on a point-to-point paradigm.

Indeed, as already described in Section 3, in the 5G network there are two options available for the architecture:

1. The Point-to-Point NG Core Architecture;

2. Service-Based NG Core Architecture (SBA);

With the SBA, there are functions that consume services and those that produce services. Any network function can offer one or more services. The framework provides the necessary functionality to authenticate the consumer and to authorize its service requests. The framework supports flexible procedures to efficiently expose and consume services. For simple service or information requests, a request-response model can be used. For any long-lived processes, the framework also supports a subscribe-notify model.

The API framework defined by ETSI ISG MEC is aligned with these principles and in fact does exactly the same for MEC applications as the SBA does for network functions and their services. The functionality needed for efficient use of the services includes registration, service discovery, availability notifications, deregistration and authentication and authorization. All this functionality is the same in both the SBA and the MEC API frameworks.

Fig. 6.1 shows how the MEC system is deployed in an integrated manner in a Service-Based Architecture 5G network.

**Figure 6.1:** 5G Service-Based Architecture and a generic MEC system architecture

In the MEC system on the right-hand side of Fig. 6.1 the MEC orchestrator is a MEC system level functional entity that, acting as an AF, can interact with the Network Exposure Function (NEF), or in some scenarios directly with the target 5G NFs. On the MEC host level it is the MEC platform that can interact with these 5G NFs, again in the role of an AF. The MEC host, i.e. the host level functional entities, are most often deployed in a data network in the 5G system. While the NEF as a Core Network function is a system level entity deployed centrally together with similar NFs, an instance of NEF can also be deployed in the edge to allow low latency, high throughput service access from a MEC host.

The distributed MEC host can accommodate, apart from MEC apps, also MEC platform services. The choice to run a service as a MEC app or as a platform service is likely to be an implementation choice and should factor in the level of sharing and authentication needed to access the service. A MEC service such as a message broker could be initially deployed as a MEC app to gain time-to-market advantage, and then become available as a MEC platform service as the technology and the business model matures.

Managing user mobility is a central function in a mobile communications system. In a 5G system it is the Access and Mobility Management Function (AMF) that handles mobility related procedures. In addition, the AMF is responsible for the termination of RAN control plane and Non-Access Stratum (NAS) procedures, protecting the integrity of signalling, management of registrations, connections and reachability, interfacing with the lawful interception function for access and mobility events, providing authentication and authorization for the access layer,

and hosting the Security Anchor Functionality (SEAF). With the SBA, the AMF provides communication and reachability services for other NFs and it also allows subscriptions to receive notifications regarding mobility events.

Similarly to the AMF, the Session Management Function (SMF) is in a key position with its large number of responsibilities. Some of the functionality provided by the SMF includes session management, IP address allocation and management, DHCP services, selection/re-selection and control of the UPF, configuring the traffic rules for the UPF, lawful interception for session management events, charging and support for roaming. As MEC services may be offered in both centralized and edge clouds, the SMF plays a critical role due to its role in selecting and controlling the UPF and configuring its rules for traffic steering. The SMF exposes service operations to allow MEC as a 5G AF to manage the PDU sessions, control the policy settings and traffic rules as well as to subscribe to notifications on session management events.

Until now the SBA has been discussed with its network functions and their roles in the 5G system. While they play an essential role in enabling the flexible integration of MEC in the next generation system, there are a few additional high-level concepts worth listing that are essential in providing high performance MEC services with an unparalleled quality of experience.

- Concurrent access to local and central Data Networks (DN) in a single PDU session;

- Selection of the User Plane Function for a PDU session close to the UE's point of attachment;

- Selection/establishment of a new UPF based on UE mobility and connectivity related events received from the SMF;

- Network Capability Exposure to allow MEC (AF) to request information about UE(s) or request actions towards UE(s);

- Possibility for MEC (AF) to influence traffic steering for a single UE or a group of UEs;

- Support for LI and Charging for MEC in the edge cloud;

- Indication about LADN availability for UEs (Local Access Data Network) for specific and local MEC services;

## 6.2.2  MEC deployment scenarios

Logically MEC hosts are deployed in the edge or central data network and it is the User Plane Function (UPF) that takes care of steering the user plane traffic towards the targeted MEC applications in the data network. The locations of the data networks and the UPF are a choice of the network operator and the network operator may choose to place the physical computing resources based on technical and business parameters such as available site facilities, supported applications and their requirements, measured or estimated user load etc. The MEC management system, orchestrating the operation of MEC hosts and applications, may decide dynamically where to deploy the MEC applications. In terms of physical deployment of MEC hosts, there are multiple options available based on various operational, performance or security related requirements. The following list outline of some of the feasible options for the physical location of MEC.



**Figure 6.2:** MEC physical deployments

1.  MEC and the local UPF collocated with the Base Station, as shown in Fig. 6.2 (a).

2.  MEC collocated with a transmission node, possibly with a local UPF, as

shown in Fig. 6.2 (b).

3. MEC and the local UPF collocated with a network aggregation point, as shown in Fig. 6.2 (c).

4. MEC collocated with the Core Network functions (i.e. in the same data centre), as shown in Fig. 6.2 (d).

# 6.3 MEC for an advanced automotive communications

Connected Vehicles and especially connected Autonomous Driving (AD) vehicles bring a whole new ecosystem with new requirements on the edge and the network architecture to support the new workloads and to satisfy the real-time service requirements. Such ecosystem includes the vehicles, the road infrastructure, the network infrastructure, and the edge.

Edge Computing based Vehicle-to-Cloud solutions enable edge cloud capabilities for different levels of autonomous driving, including Highly Autonomous Driving (HAD) and Fully Autonomous Driving (FAD) through providing different services for the driving process (e.g., High Definition real-time Maps, real-time traffic monitoring and alerts, and richer passengers experience), supporting vehicles on roads to drive cooperatively and to be aware of road hazards, and providing better user experience and trust to drivers and passengers.

## 6.3.1 MEC as solution for for the AD applications

The Multi-access Edge Computing (MEC) is taking a predominant role in the automotive context as a standardized solution for Edge Computing, especially important from automotive stakeholders' point of view. In particular, from a standardization perspective, some use cases targeting fully connected cars (i.e. FAD with the maximum level of automation) have challenging requirements that may be fulfilled only with the introduction of the MEC in the next generation 5G networks.

The expansion from Cloud to MEC for connected AD services is driven by both the need to have more processing power closer to the vehicles to guarantee the required latency and the need to have reduced network churn with continuous access to the Cloud. MEC is addressing this paradigm shift by aiming to offer a

different services environment and cloud-computing capabilities within the roads infrastructure and the access network infrastructure in close proximity to vehicles and Road Side Units (RSUs).

MEC can also benefit from new functionality to better handle the big volume of data coming from vehicles and road side units and to dynamically allocate CPU and Acceleration resources based on the services' needs (e.g., computer vision Vs. video streaming Vs. data aggregation).

The information exchanged between vehicles, infrastructure, pedestrians, and network using V2X technology is enabling a multitude of new and exciting applications. Exploitation of the edge processing power and its ability to intelligently process the information can add value to it, and to provide useful low latency service experiences.

MEC provides application and content providers with cloud computing capabilities and IT service environment at the very edge of the mobile network. This environment is characterized by the proximity, often in both physical and logical sense, to the clients, enabling very low latency between the client and the server applications, high bandwidth for the application traffic, and near real-time access of the applications to context-rich information, e.g. related to device locations and local radio network conditions. These qualities of MEC ensure an unparalleled quality of experience with highly contextualized service experience and efficient utilization of radio and network resources.

In addition, MEC can be deployed within the Mobile Network Operator's infrastructure together with the management and orchestration, security, privacy and subscriber management frameworks already in place. This way the environment for edge applications is secure and well managed, making it more suitable critical applications as well as for applications with high business value.

## 6.3.2 MEC in the Automotive use cases

The 5G Automotive Association (5GAA) categorizes a comprehensive list of connected vehicle applications, categorized in four main groups of use cases [50]:

1. Safety,

2. Convenience,

3. Advanced Driving Assistance,

| V2X Groups | Use Cases | Description | MEC Relevance |
|---|---|---|---|
| **Safety** | Intersection Movement Assist | Warn driver of collision risk through an intersection | High |
| **Convenience** | Software Updates | Deliver and Manage Automotive Software Updates | Mid |
| **Advanced Driving Assistance** | Real-Time Situational Awareness & High Definition Maps | Alert driver of Host Vehicle (HV) moving forward of hazard road conditions in front | High |
| | See-Through | Driver of HV that signals an intention to pass a Remote Vehicle (RV) using the oncoming traffic lane is provided a video stream showing the view in front of the RV. | High |
| | Cooperative Lane Change (CLC) of Automated Vehicles | Driver of HV signals an intention to change the lane with at least one Remote Vehicle (RV) in the target lane in the vicinity of the HV | High |
| **VRU** | Vulnerable Road User Discovery | Detects and Warns drivers of VRUs in the vicinity | High |

**Table 6.1:** V2X Groups & MEC Relevance

4. Vulnerable Road User (VRU).

The V2X use case group "Safety" are designed to reduce the frequency and severity of accidents. The United States' National Highway Traffic Safety Administration has compiled some extensive research and statistics regarding both vehicle-to-vehicle [51] and vehicle-to-pedestrian [52] crash scenarios. Such use case group includes several different types of use cases to support road safety using the vehicle-to-infrastructure (V2I) communication in addition to the vehicle-to-vehicle (V2V). For some use cases, MEC systems could provide a support for Real-time data analysis, data fusion and reduced ingress bandwidth with respect to the remote cloud.

The V2X use case group "Convenience" provides time-saving services to manage data and the health of the vehicle. This group of V2X use cases requires

cost-effective communication to be enabled between the vehicles and the backend server (e.g., car OEM's server). Software Over the Air (OTA) updates and other telematics use cases are typically included in this group.

The V2X use case group "Advanced Driving Assistance" are focused on improving traffic flow, traffic signal timing, routing, variable speed limits, weather alerts, etc. This group of use cases collects the most challenging requirements for V2X (from a MEC perspective). It can require distribution of a relatively large amount of data with high reliability and low latency in parallel. Additionally, the advanced driving use cases would benefit from predictive reliability. This means that vehicles moving along should have the possibility to receive a prediction of the network availability ahead of them to allow preparations accordingly.

The V2X use case group "VRU" supports a safe interaction between vehicles and pedestrians, motorcycles, bicycles or any other non-vehicle road user.

Table 6.1 summarizes the V2X use case groups, their applications and the relevance for the MEC.

The use cases listed in Tab. 6.1 identify several challenges from a technical standpoint. MEC solutions can contribute to the effective realization of many V2X use cases of interest, especially when performance and system requirements are challenging, in terms of low delay, QoS management and prediction, deployment flexibility and access to local context-rich information. In addition to that, the two challenges, big data processing and storage, and multi-operator support, are particularly important to be solved for a successful implementation of automotive use cases.

Therefore to demonstrate the benefits deriving by the utilization of the MEC in the automotive context two applications belonging to the Advanced Driving Assistance V2X Groups have been developed and tested in the ARNO-5G testbed. Such applications are deeply described in the following sections.

## 6.4 A MEC-Based VRU Warning System

As already described in Section 6.3, exploiting the MEC architecture, the development of applications aimed at safety of vulnerable road users such as pedestrians, motorcyclists, public transport users, non-motorized vehicles and bicyclists becomes possible. In this context, a system addressing the use cases "Warning to Pedestrian against Pedestrian Collision" described in [53] is presented in the following of this section.

## 6.4.1 System Architecture

The proposed application architecture is based on VRU/V2N communication and, in particular, it focuses on the scenario shown in Fig. 6.3 where a VRU and a vehicle are connected to the same evolved Node B (eNB). The V2X messages, e.g. Cooperative Awareness Message (CAM), are used by the V2P application and sent to a server deployed at the network edge or in the cloud. A block diagram of the proposed architecture is shown in Fig. 6.4.



**Figure 6.3:** MEC-Based VRU Warning System Architecture

The CAM client deployed through an Android application, described in the following, sends periodically CAMs to the CAM server that in turn by means of the Listening CAM module parses each messages and stores the information related to the sending client. Then the distance between the VRUs (e.g., the pedestrian) and every nearby entities is computed by using Federal Communications Commission approved formula of ellipsoidal Earth projection by the GEO-Computation Engine. Finally, the Logic Sending Alerts sends an alert to the involved road entities if such distance is less than a fixed threshold. To take trace of the evolving situation the Client for Demo Viewer sends to the Web Socket GUI the needed information for a web browser visualisation.



**Figure 6.4:** MEC-Based VRU Warning System Block Diagram

## 6.4.2 CAM Client: the VRU app

The pedestrian awareness Android application is the starting point of the entire implementation. The application monitors every five seconds the geographical location of the VRU and send a CAM to the CAM server. After starting the application, through a drop-down menu as shown in Fig. 6.5, the VRU can set the CAM server IP Address where send the CAM and start or stop its position acquisition. In the same drop-down menu are also shown the last acquired position (e.g Latitude and Longitude), the speed and the direction of the VRU. The application is developed using Android Studio 3.1 and Google Play services SDK tool. The Google Play services SDK tool permits to use the newest APIs in order to get the VRU position, create the map on the smartphone screen and put a marker on the acquired geolocation data. The Google Play services tool permits also an easy interaction with all the other services used by the application, without worrying about device support.



**Figure 6.5:** MEC-Based VRU App

The application is capable to understand the pedestrian position, exploiting the on-board GPS transceiver of every Android smartphone. Indeed, the latitude and longitude coordinates of a new position are acquired through the GPS Android APIs `GPS.getLatitude` and `GPS.getLongitude`.

Once the coordinates are acquired, a red marker is set on the new acquired position and depicted on the map exploiting the `googleMap.addMarker` API. Furthermore, through the googleMap.moveCamera API, the application moves the red marker at the center of the smartphone screen. A trace of the last positions

is taken with a blue marker. All the markers show the latitude and longitude coordinates if clicked.

Finally, the application, through both Wi-Fi and LTE network, is able to connect and send the position data remotely to the CAM server.

The pedestrian awareness Android application has been tested in a Huawei P8 Lite smartphone hosting Android 7.0 OS.

## 6.4.3 Collision Avoidance Messages (CAMs)

In the Intelligent Transport System (ITS) the messages exchanged between the stations (ITS-Ss) composing it, are named Cooperative Awareness Messages. They are used to create and maintain awareness of each other and to support cooperative performance of the road users.

A CAM contains information of the originating ITS-S, and it depends on the type of involved ITS-S. When an ITS-S receives a CAM, it becames aware of the presence, type, status of the originating ITS-S. Therefore, the receiving ITS-S could use this information for comparing the status of the originating ITS-S with its own status and compute the collision risk with the originating ITS-S. CAM are generated periodically, with a frequency of 1- 10 Hz by the Cooperative Awareness (CA) basic service in the originating ITS-S. Such frequency takes into account the change of own ITS-Ss status, e.g. the change of the position or speed. CAM messages can be up to 800 bytes in length and they have to experiment a max latency of 100 ms.

In the proposed scenario, CAMs are generated according the library [54]. Each CAM contains:

- Station ID;

- Message ID;

- Position: latitude and longitude expressed in micro-degrees and altitude;

- Speed, expressed in m/s;

- Drive Direction (N,N-E,E, S-E, S, SW, N-W).

A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is

following:A CAM sample is following:A CAM sample is following:A CAM sample is following:A CAM sample is following:

**Listing 6.1:** CAM example

```
CAM(Header(protocolVersion 1, messageID 0, stationID 0),
   CoopAwareness(generationDeltaTime=100, CamParameters(Basic:
      ReferencePosition(43718285, 10424889, PosConfidenceEllipse
      (4095, 4095, 3601), Altitude(800001, unavailable)), Speed(
      speedValue(0), speedConfidence(null)), driveDirection(SOUTH)))
      )
```

## 6.4.4 CAM Server

A multi-threaded software daemon has been developed to act as CAM server. The working algorithm is depicted in Fig. 6.6 and it can be described as follow:



**Figure 6.6:** CAM Server Flow Chart

1. The server which listens for incoming connections on a designed port;

2. Vehicles connect to the server, and then send a CAM message;

3. The server parses the CAM message and verifies that the syntax is correct;

4. Then it stores the information into a data structures until it is useful, to let the software to access it easily. VRU messages and vehicles messages are stacked into different arrays;

5. For each vehicle message is computed the distance between the vehicle and the VRUs. The distance computation is performed by using FCC-approved formula of ellipsoidal Earth projection [55], that reduces the complexity of the computation while its approximation is valid for usual distance between vehicles and VRUs.

6. If the distance between VRU and a vehicle is lower than a given threshold, the server sends an alert Decentralized Environmental Notification Message (DENM) to the VRU. The threshold is dynamically computed according the VRU and vehicle speed.

## 6.4.5 VRU Viewer

To visualize the evolving situation the CAM server has been extended with an optional module that sends information to dynamic browser page, named demoViewer. CAM server sends information by using UDP to the nodejs gateway that is able to create a persistent connection with the browser.



**Figure 6.7:** MEC-Based VRU App Viewer

The web page dynamically shows a map centered on eNB coordinates, the vehicles and the VRUs. The web page contains the JavaScript that setups the page

**Figure 6.8:** Experienced CAM Latency from VRU to CAM server.

loading the map from openstreet servers, then it connects to the nodejs gateway to continuously receive information from the CAM server about VRU, vehicles, and threshold, and finally it updates VRU and vehicles positions on the maps according the received coordinates. If a proximity alarm information is detected, i.e. when the distance between the VRU and the vehicles is less than the threshold, it shows an popup alert over VRU head, containing a warning message and the distance information.

The Fig. 6.7 illustrates an example of the demoViewer web page: here is shown the 5G antenna, the VRU (yellow man), and three vehicles as oriented green pointers.

## 6.4.6 VRU Warning System Performance evaluation

Scenario with one VRU and multiple emulated vehicles has been setup. The VRU position was fixed while vehicles were moving according a predefined path at a variable speed between 10 and 30 km/h. Emulated vehicle CAM messages were sent every 1s (as per CAM standard). In a first set of measurements, the CAM message time between the mobile device and CAM server located in the ARNO-5G testbed is considered. As depicted in Fig. 6.8, when a WiFi connection is considered, the average Round Trip Time is about 9.6 ms (min. 2.7 ms - max. 29.2 ms), while if the mobile device is attached to the 5G OAI network, the average RTT is

about 25.3 ms (min. 15 ms - max. 31.9 ms).

Then the CAM server is moved on a cloud server, hosted at OVH provider. The message time between device and CAM server in such configuration increased to 37.3 ms (min. 25.3 - max. 46.5 ms) as shown in Fig. 6.8.

Accordingly to the obtained results, the factor that most contributes to latency between the mobile device and the CAM server is depending by the wireless network access time.

## 6.5 Orchestrating Heterogeneous MEC-based Applications

As already described in the previous sections of this chapter, the applications grouped in the ADA use case (refer to Tab. 6.1 for more details) are designed to improve traffic flow, traffic signal timing, routing, variable speed limits and weather alerts. This group of applications is characterized by the most challenging requirements in terms of latency and requested bandwidth. Indeed, they require the distribution of an often large amount of data with high reliability and low latency.

At the same time, video streaming applications are gaining popularity and are becoming one of the major Internet services for mobile consumers. Thus, in-vehicle infotainment services are expected to stand out among the most popular connected vehicle services.

MEC can furthermore facilitate the cooperation of telecom operators and the automotive vertical industry and play a fundamental role in this direction: MEC provides an execution environment for the deployment of third party applications at the mobile edge with a significant level of network awareness via standardized interfaces, such as the Radio Network Information Service (RNIS) API [42]. This has the advantage of not only serving end users from edge servers, thus minimizing latency and reducing core network traffic, but also enabling the vertical service provider to perform network-aware optimizations exploiting real-time radio network information, such as network load and per-user channel quality indications.

This section presents the design, implementation and evaluation of a Connected Vehicle Service Orchestrator (CVSO), which aims to optimize the QoE of a video streaming service for infotainment, while guaranteeing the requested capacity to coexisting ADA services. The orchestrator is deployed as a MEC application, as

is also the case for the two considered services. The orchestration algorithm involves the assignment of the appropriate video quality to each user in order to maximize the overall viewing experience; this problem is formulated as an Integer Linear Program (ILP). The orchestrator consists of a Video Streaming Controller (VSC) that implements and solves the ILP problem and a Radio Link Measurement (RLM) component, which monitors the network status by accessing the RNIS and feeds such information to the ILP model. Based on the obtained results, the VSC sets the allowable video streaming bitrate, and thus quality, for each mobile user that the Video Streaming Server (VSS) application delivers. We should note that video delivery is implemented using Dynamic Adaptive Streaming over HTTP (DASH) [56] technologies. This has the advantage of wide compatibility with standard HTTP servers and video players.

An evaluation of the proposed scheme using the 5G-ARNO testbed over a particular version of OpenAirInterface (OAI) platform, which also includes a standards-compliant MEC platform implementation is performed. The results show that the CVSO brings significant improvements in terms of user experience compared to standard receiver-driven DASH video adaptation mechanisms and to non-adaptive streaming.

Finally, an experimental performance assessment of our scheme in terms of the execution time to solve the ILP problem validates its applicability for MEC deployment, where compute resources are typically scarcer. It is also demonstrated that the time to solve the ILP is in the order of seconds using a standard solver even for very large numbers of users, which adds to our system's feasibility. The presented experimental results we on the algorithm's running time as a function of the CPU resources assigned to it can provide insight on how compute resources should be allocated to the CVSO, in order to maintain specific response time goals; this is important in an environment where low latency matters.

## 6.5.1 CVSO Architecture and Components

This section details the Connected Vehicle Service Orchestrator (CVSO) architecture and its MEC-based components, namely the Video Streaming Controller (VSC) and the Radio Link Measurement (RLM) module. In addition, the Video Streaming Server (VSS) application is described.

Fig. 6.9 shows the proposed orchestrator architecture. The MEC Platform hosts the two orchestrator components implemented as MEC applications: the VSC and

**Figure 6.9:** CVSO Architecture

the RLM. The VSS is also implemented as a MEC application. The ADA component (road safety application) is depicted in red in Fig. 6.9. Its capacity requirements are known to the CVSO, and are taken into account by subtracting the amount of capacity that it necessitates from the available radio link capacity for each connected vehicle. The internal workings of this component are outside the scope of this work.

## Video Streaming Control and RAN monitoring

The *VSC* takes the orchestration decisions. It consists of three main functions:

(i) It periodically builds an ILP problem to assign video representations/qualities to each connected user accessing the video service, taking into account the Channel Quality Indicator (CQI) values received from the *RLM*.

(ii) It solves the above optimization problem.

(iii) According to the solution of the optimization problem, it selects the appropriate video representation for each user and communicates it to the VSS as described in the following section.

The purpose of the *RLM* is to maintain a view of the status of the RAN conditions. In particular, it accesses the RNIS over its standardized interface in order to have an up to date view of the number of users per eNodeB in the area that it

manages and their Channel Quality Indication (CQI) values. This information is critical in order to estimate each user's link capacity, and, in turn, the video bitrate that it can sustain.

It should be noted that each edge data center and the respective MEC platform may correspond to a region containing multiple mobile network cells. In this scenario, each CVSO would orchestrate the video delivery and ADA services for all the users/vehicles within these cells. For simplicity, we can consider that there is a single MEC Platform, and thus, CVSO instance, per network cell.

The CQI values are normally reported by the UEs to the eNodeB via standard LTE procedures. In our MEC platform implementation, which is tailored to OAI, we use the FlexRan architecture and protocol [57] to extract them (as well as a wealth of other RAN-level information) from the eNodeB and make them available to the RNIS, so that MEC applications such as the RLM can consume them. For more details on the considered MEC platform and RNIS implementation, the reader is referred to [58].

The *VSC* receives CQI updates from the *RLM* and uses this information, among others, as ILP input.

## MEC-assisted video streaming service

The *VSS* contains the video available for streaming and the information describing the contents. In our proposed scheme the *VSS* is also deployed as a MEC application instance and plays the role of an edge video cache. This has the advantage of serving users from a nearby location, thus reducing startup latencies and saving on backhaul network resources. Our adopted technology for video streaming is Dynamic Adaptive Streaming over HTTP (DASH). DASH is a streaming technique that enables high quality streaming of media content over the Internet delivered from conventional HTTP web servers. MPEG-DASH works by breaking the content into a sequence of small video segments, each containing a short duration of video (typically a few seconds). Information about the available representations needed by clients is stored in the Media Presentation Description (MPD) file [59]. The MPD file structure follows a hierarchical data model containing one or more *periods*. Each period contains video characteristics like available bitrates, resolutions, codecs, segment IDs, etc.

Normally, to play the content, an MPEG-DASH client first obtains the MPD using the HTTP protocol. Analyzing the MPD file, the client learns about all the

content characteristics. Then, it selects according to the network conditions the appropriate representation among the available listed in the MPD. Finally, it starts streaming the content by fetching the segments related to the selected representation using HTTP GET requests. After appropriate buffering to deal with network throughput variations, the client continues fetching the subsequent segments of the considered representations. Meanwhile, the client may monitor the available network bandwidth and its fluctuations and appropriately adapt to it by fetching segments of a different representation characterized by a lower or higher bitrate, as listed in the MPD to match the current network conditions.

One disadvantage of pure receiver-driven approaches for bitrate adaptation is the fact that the video client takes adaptation decisions based on bandwidth estimates only, lacking accurate knowledge of the actual network conditions. This is one of the issues that the MPEG Server And Network Assisted DASH (SAND) standard [60, 61] aims to address. SAND involves the exchange of quality-related information between SAND-aware clients and servers and between DASH-Aware Network Elements (DANE) for improved adaptation decisions. SAND, however, requires specific extensions to DASH clients to support it.

The proposed approach to offer server assistance in the video delivery process is different, since one of our targets was for our server-side mechanisms to be transparent to the clients and to support default DASH players without any modifications. In the proposed scheme, since the VSS is deployed at the MEC, it can take advantage of the availability of RAN-level information coming from the RNIS and adapt the contents of the MPD files requested dynamically and on a per-user basis, according to the current link capacities of the involved clients. We have implemented this functionality as follows. The VSS operates an HTTP front-end proxy which receives user requests for MPD files. This proxy also has an internal REST HTTP interface where the VSC posts updates about the outcome of the video representation assignment algorithm. In particular, each time the VSC is executed, it updates the VSS front-end with ¡user IP address-video representation¿ pairs. When the VSS front-end receives a request for an MPD file, it fetches it from the VSS, modifies it so that it contains only the appropriate video representation, and serves it back to the user. This procedure is transparent to the user, which proceeds with downloading the video chunks that correspond to the representation selected for it by the VSC.

To force the clients request periodically an updated MPD that reflects the current network state, we use a mechanism that is specified in MPEG-DASH, making

our solution standards compliant: we exploit the `minimumUpdatePeriod` MPD attribute. When this attribute is set in the MPD file, the client re-fetches a fresh version of it when the timer defined in the `minimumUpdatePeriod` expires.

## 6.5.2 An ILP Formulation for QoE-aware video representation selection

In the considered application scenario, the video service provider makes videos available in a distinct set of representations, each with its own bitrate which maps to a specific QoE value. It should be noted that for the same codec settings, bitrate, and resolution, and even if we ignore the effects of other impairments (buffering events, display qualities, etc.), different videos might come with different QoE since the content itself (e.g., the amount of motion in the video) affects the relationship between video quality and bitrate. In this study, such effects are not taken into account.

The purpose of the VSC is to select the most appropriate video representation for each user, taking into account the available network capacity, load (i.e., number of users sharing radio resources), and the individual link characteristics (channel quality), in order to maximize the overall user experience expressed as the sum of the user QoE values. This problem has been formulated as an ILP with the following parameters:

- $n$: number of mobile users connected to the *VSS*;

- $m$: number of video representations available in the *VSS*;

- $b_j$: video bitrate of representation $j$;

- $E_j$: QoE corresponding to video representation $j$;

- $R^i_{1PRB}$: bitrate achievable by user $i$ in one Physical Resource Block (PRB). It is extracted from *Table 7.1.7.1-1* of [62] and it is a function of the CQI value of the $i^{th}$ user;

- $R^i_{MAX}$: maximum bitrate achievable by user $i$ accordingly to its CQI value. This value is computed by assuming that user $i$ is assigned all the PRBs available for the video streaming service and it is extracted from *Table 7.1.7.1-1* of [62] according to the CQI value of the $i^{th}$ user.

- $PRB^{VSS}$: portion of the overall number of wireless channel PRBs dedicated to the video streaming service;

The binary variable $x_{i,j}$ indicates that mobile user $i$ is assigned video representation $j$.

In our problem formulation, QoE is expressed in terms of the Mean Opinion Score (MOS). The MOS is defined as the expected rating that a panel of users would give to the quality of the transmitted video in the 1-5 scale, where 1 represents the poorest quality and 5 an excellent one [63]. QoE can be estimated from video quality parameters, such as interruption statistics and encoding parameters, which can be translated to MOS as we explain in next section. The estimated MOS that corresponds to each representation ($E_j$ values), which is related to its bitrate and, in turn, to the specific encoder parameters used to produce it, is reasonably assumed to be known to the VSC.

The problem is therefore modelled as the following binary ILP.

**Objective Function**:

$$Max \sum_{i=1}^{n} \sum_{j=1}^{m} x_{i,j} E_j \tag{6.1}$$

**Constraints:**

$$\sum_{i=1}^{n} \sum_{j=1}^{m} \frac{b_j}{R_{1PRB}} \leq PRB^{VSS} \tag{6.2}$$

$$x_{i,j} b_j \leq R_{MAX}^i, \forall i \in [1, n] \tag{6.3}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{m} x_{i,j} \leq 1, \forall i \in [1, n] \tag{6.4}$$

Constraint (6.2 ensures that the sum of all the assigned PRBs to all the mobile users requesting the video streaming service is less than the bandwidth dedicated to the video streaming service. In this way the assigned bandwidth to the video streaming is not exceed and the bandwidth required by the ADA services is assured and available.

Constraint (6.3) ensures that a video representation is not assigned to user $i$ if its bitrate cannot be accommodated by the user. Finally, constraint (6.4) guarantees that a user is assigned at most one video representation.

## 6.5.3 Testbed description

The testbed shown in Fig. 6.10 is utilized for a preliminary experimental evaluation.



**Figure 6.10:** ARNO-5G Configuration for CVSO performance evaluation

A complete MEC system tailored to OAI that complies with the ETSI MEC specification has been implemented. It also includes a fully-fledged RNIS. The considered testbed includes the appropriate extensions to interface the 5G OAI network with the MEC platform. The OAI-based Evolved Packet Core (EPC) contains the implementation of the following core network elements: the Serving Gateway (S-GW), the PDN Gateway (P-GW), the Mobile Management Entity (MME) and the Home Subscriber Server (HSS). These are deployed as a set of virtual machines on top of the kvm hypervisor. At the S/P-GW level we have followed the Control-User Plane Separation principle [64], and use it to implement traffic steering towards MEC applications. In particular, the S/P-GW control and data plane are split, and user data plane traffic is handled by an OpenVSwitch (OVS). When specific MEC applications are deployed to which traffic needs to be offloaded (as is the case, e.g., for the VSS), the appropriate traffic redirection rules are installed by the MEC platform to OVS and the latter steers traffic appropriately. Further details on our MEC platform are not within the scope of this paper.

The OAI eNB is a standards-compliant implementation of a subset of LTE Release 10. The Ettus B210 USRP device acts as radio front-end, and it is connected to PC 8 running the OAI eNB software.

Two Huawei EE372 dongles, each connected to a laptop, equipped with a preconfigured SIM card able to connect to the OAI mobile network act as User Equipment (UE). The UEs are connected to the same eNB and in order to have significantly different CQI values are located at a different distance from the eNB. A first UE (i.e., $UE_1$) is located at a distance $d_1$ while a second UE (i.e. $UE_2$) is placed at a distance $d_2$.

The MEC platform hosts the *CVSO* and its MEC-based components, namely the

*VSC* and the *RLM*, as well as the *VSS*. The RLM consumes CQI reports from the RNIS and makes them available to the VSC in order to run the video representation assignment algorithm. The latter communicates with the VSS front-end to notify it of up-to-date video representation assignments over the latter's internal interface via HTTP.

## 6.5.4 Experimental Methodology

In the performed experiments, different representations of the same test video (the 10m54s Blender Foundation's "Elephants Dream" [65]) are used. Such representations are stored in the *VSS* and are made available for streaming. Using the *ffmpeg* tool [66], the test clip is H.264/AVC-encoded in six files with different bitrates approximately equal to $256\ kbps$, $512\ kbps$, $1024\ kbps$, $2048\ kbps$, $4096\ kbps$, and $8192\ kbps$, respectively, with a $1280 \times 720$ resolution.

The different representations are then prepared for DASH delivery using GPAC's MP4Box utility [67]. Each segment has a duration of $2s$. According to Lederer [68], this is an appropriate segment duration for an environment with low delay (due to MEC-based delivery) but high bandwidth fluctuations (due to the potentially varying channel conditions of mobile users), helping clients adapt faster to bandwidth changes and limit stalls.

The `minimumUpdatePeriod` MPD file attribute is set so that the mobile users request an updated version of the MPD file every $10\,s$.

The proposed *CVSO* is then evaluated in terms of the achieved QoE when the "Elephants Dream" movie is streamed simultaneously to both UEs. QoE is principally affected by the video encoding quality, determined by the Quantization Parameter (QP) used for encoding and the Playout Interruptions (PI), caused by buffering delays due to insufficient bandwidth. We make the assumption that the effect of the impairments related to interruptions and picture quality (QP) on QoE is additive. The theory that quality impairment factors have an additive effect on the psychological scale is due to Allnatt [69]. This is a typical assumption for QoE assessment, both for video [70, 71] and VoIP services [72]. In order to measure these impairments, we use a combination of existing tools as in [73]. Regarding the effects interruptions ($I_{PI}$), we use the Pseudo-Subjective Quality Assessment (PSQA) approach [74], and in particular the PSQA tool of Singh et al. [75], configuring it to ignore QP-related impairments. The PSQA methodology involves training a Random Neural Network (RNN) using data from subjective tests. The

trained RNN classifier can then be applied to calculate the expected MOS for specific values of the input parameters.

On the other hand, we apply the QoE measurement tool developed in the context of the VIPEER project [76] for measuring picture quality impairments ($I_{QP}$), which ignores the effects of interruptions and only considers the average value of the QP for each QoE measurement window (16 s of video, in our case). Subtracting the $I_{QP}$ and $I_{PI}$ values from the maximum possible MOS corresponding to excellent quality (i.e., 5), we get the MOS value for the specific measurement window:

$$MOS = 5 - I_{QP} - I_{PI}. \tag{6.5}$$

In order to be able to measure and export the QP and PI information necessary for our QoE calculations, we have appropriately extended the open-source MP4Client video player and the libx264 video encoding library and installed it at the UE hosts.

## 6.5.5 Experimental Results

In the experimental measurements, the following four different scenarios are considered.

(i) In the first scenario the *CVSO* is exploited to optimize the video bitrate. The experimental results obtained in such scenario and labelled as "RAN-aware" are depicted by means of green boxes curves in Fig. 6.11 and Fig. 6.12.

(ii) In the second scenario the video bitrate is fixed to a low value for all users. Such experimental results labelled as "Fixed [low quality]" are depicted by means of blue crosses curves in Fig. 6.11 and Fig. 6.12.

(iii) In the third scenario the video bitrate is fixed to a medium value for all users. Such experimental results labelled as "Fixed [medium quality]" are depicted by means of pink dots curves in Fig. 6.11 and Fig. 6.12.

(iv) In the last scenario, the default video bitrate adaptation is used. The experimental results obtained in such scenario labelled as "Receiver-driven" are depicted by means of orange stars curves in Fig. 6.11 and Fig. 6.12.

In scenario (i), the *VSC* being aware of the CQI values of each mobile user accessing the video service is able to optimize the MOS of $UE_1$ for the entire video

**Figure 6.11:** QoE for $UE_1$ (CQI=15)

duration obtaining the higher MOS values (equal to 4) as shown in Fig. 6.11. Regarding $UE_2$, as shown in Fig. 6.12, the performance is naturally lower due to a lower CQI value.

In the second and third scenario, there is no awareness of the network conditions and the video streaming bitrate is fixed. In the second scenario, a low video quality resolution corresponding to a low video streaming bitrate is imposed to satisfy also the video streaming demand of the most disadvantaged UE (i.e., $UE_2$). Due to the low video streaming bitrate adopted in this scenario, the MOS experienced by all the mobile users is less then the one experienced in the first scenario where the video streaming bitrate is optimized. Then in the third scenario, to increase at least the MOS of the most advantaged UE the video quality resolution is improved selecting a higher video streaming bitrate. The MOS values of both UEs are increased as depicted in Fig. 6.11 and Fig. 6.12, but, in particular for $UE_1$, it is lower then the MOS obtained in scenario (i).

In the fourth scenario, the adaptation of the video bitrate is left to the DASH client. The mobile user that start the streaming first (e.g., $UE_2$) requests all the available signal bandwidth to the eNB to maximize its own MOS. The other mobile user (e.g., $UE_1$), not having enough resources, experiences a very low MOS, as depicted in Fig. 6.11 and Fig. 6.12. It shall be noted that for both UEs after a few seconds the video streaming crashes.

**Figure 6.12:** QoE for $UE_2$ (CQI=10)

Thus, the receiver-driven curves in both Fig. 6.11 and Fig. 6.12 are truncated. In addition, the MOS obtained by the $UE_1$ in this scenario is much lower than the MOS obtained when the *CVSO* is exploited.

Finally the average of MOS of the $UE_1$ and $UE_2$ for all the above mentioned cases is calculated. As shown in Fig. 6.13, the the case for using the *CVSO* scores consistently better.

## 6.5.6 ILP solution performance

Turning our attention to the performance of the *CVSO* with respect to solving the ILP for video representation assignment, we perform a set of experiments to measure its execution time.

The ILP solution time is defined as the time needed by the *VSC* module to derive a solution to the ILP, once the latter has been configured in the solver. The overall solution time instead is defined as the sum of the time needed by the *VSC* module for creating the constraints of the ILP problem, solving the ILP problem and for setting the video streaming bitrate according to the ILP solution results.

The model presented in Section 6.5.2 is solved using the IBM ILOG CPLEX Optimizer [77].

In the performed measurements, the number of mobile users and the number

**Figure 6.13:** Average QoE Comparison

of CPUs assigned to the *VSC* are increased. The number of mobile users is varied from 10 to 10000, while the number of used CPUs is varied from one to four. Each measurement is repeated 50 times. Results are reported with 95% confidence intervals.

Fig. 6.14 and Fig. 6.15 show the time needed for the CPLEX problem solution as function of the connected mobile users for all users and for up to 100 users, respectively.

The results show that the required time for solving the ILP (i.e. CPLEX Solving Time) is consistently less than the `minimumUpdatePeriod`. Therefore the proposed *CVSO* is able to adapt the video streaming bitrate of each user at the network conditions before any MPD file modification performed by the *VSC*. This conclusion is even more valid if the number of CPUs assigned to the problem solution is increased.

The overall time shown in Fig. 6.16 grows despite the increasing number of CPUs assigned to the *VSC*. This is because the constraints are set anew every time CPLEX solves the problem and the parallelism benefits deriving from the assignment of more than one CPU are not exploited in this phase of the solution process. However please note that when 10000 mobile users are connected to the same eNB, very overestimated condition, the overall solving time is less than one minutes. Indeed as stated in [78], a generous and practical evaluation of the num-

**Figure 6.14:** CPLEX Solving Time



**Figure 6.15:** CPLEX Solving Time for the first 100 mobile users

**Figure 6.16:** Overall Solving Time

ber of simultaneously connected users is up to 100.

## 6.6 Conclusions

ETSI has published the baseline MEC standards to allow a standards-based environment for cloud applications in the network edge. Application enablement API and the MEC service APIs are the essential components in the MEC specification for a unified, standards-based environment for context-aware cloud applications.

New MEC service APIs are also being developed for specific industry applications such as V2X to allow MEC better serve and add value to these applications.

The 3GPP 5G system specification of Rel-15 includes native enablers for edge computing. Section 6.2 has illustrated the potential of these enablers for an integrated MEC deployment in 5G networks. The key components of this integration are the ability of MEC, as a 5G AF, to interact with the 5G system to influence the routing of the edge applications' traffic and the ability to receive notifications of relevant events, such as mobility events, in the 5G system for improved MEC deployment efficiency and end user experience. Moreover, the versatility of the 3GPP service exposure and API frameworks in principle also allows MEC to provide services to the 5G system.

Section 6.3 provides an overview of automotive use cases, as introduced by 5GAA, and shows how MEC can be considered as a key technology supporting multiple services for connected AD vehicles.

Moreover, this section draws attention to the value of MEC as a standardized solution for Edge Computing, especially important from automotive stakeholders' point of view while also serving other vertical market segments.

In fact, a great value is associated with the standardization of edge computing technology, as open standards are the way to open the market and to ensure interoperability. In particular, from a standardization perspective, some use cases targeting fully connected cars will require the fulfillment of challenging requirements, possible only with the introduction of 5G networks.

In this perspective, Section 6.4 describes a MEC-Based VRU Warning System while Section 6.5 proposes a MEC-Based orchestrator to perform a needed meticulous orchestration of different MEC services such as Advanced Driving Assistance and in-vehicle infotainment services. The orchestrator proposed in Section 6.5 aims to optimize the QoE of a videostreaming service for infotainment, while guaranteeing the requested capacity to coexisting ADA services.

In particular, Section 6.4 presents a complete standards-compliant design and implementation of a Vulnerable Road User Warning system. This system aims to provide an exchange of information between the road users (e.g. pedestrian, cyclist, vehicles) about the presence of nearby entities in case of dangerous situations: warnings are provided to VRUs to avoid collision with the moving vehicle and vice versa. The proposed architecture is composed by a user-side Android application that periodically sends messages containing VRU position, speed and orientation, and by a MEC-based application, named CAM server, able to process the geographical location of the VRU, predict a possible collision and, if needed, warn the involved road entities sending alert messages. A viewer showing a map with VRU and moving vehicles has been developed to help to visualize the involved entities and as demonstrator of a viewer that could be installed both on vehicles and VRUs. Moreover, in this section a preliminary performance evaluation about end-to-end latency between VRUs application and the CAM server is presented.

Taking advantage of the recent advances in the area of Multi-access Edge Computing, the design and the implementation of a *Connected Vehicle Service Orchestrator (CVSO)* that manages heterogeneous automotive applications at the mobile edge is presented in Section 6.5. The presented system delivers QoE-optimized infotainment video services coexisting with Advanced Driving Assistance (ADA)

applications, and does so in a fully standards-based way: video is delivered using DASH technologies, which are widely supported, and the proposed server-side video quality optimizations take place in a transparent way to the clients.

At the same time, the proposed CSVO is deployed on top of a standards-compliant MEC platform featuring a Radio Network Information Service, complying with recent ETSI standards that we have implemented. Contrary to typical receiver-driven video adaptation mechanisms, it is our CVSO that decides on the optimal video quality per user, based on the latter's channel quality characteristics which are available via the RNIS. This involves solving an ILP.

Via testbed experiments over a MEC-capable LTE network, it has been shown the proposed server-assisted video adaptation scheme to improve on user experience, while also showing that the ILP solution can be derived in reasonable time even for large numbers of users, thus verifying the proposed system's feasibility and suitability for MEC deployment.

# 7 Reliability in 5G networks

The 5G network architecture, will be heavily based on virtual network functions (VNFs). Network Function Virtualisation (NFV) enables an easy introduction of new network services by adding dynamic programmability to network devices (e.g., routers, switches, and applications servers) that, in turn, empowers fast, flexible, and dynamic deployment of new network and management services. Moreover, NFV also enables network slicing by providing multiple instances of the same network function. In this context, the dynamic service chaining allows the delivery of a new breed of applications (e.g., cloud robotics, smart cities) by dynamically selecting and composing computational and network services deployed as virtual functions (VFs) in distributed micro-clouds located at the network Edge closer to the users.

The exploitation of NFV is foreseen also in the NG Core [79] and the NG RAN technology. In the NG Core for example, the different network functions (e.g., Access and Mobility Function (AMF), Session Management Function (SMF), Policy Control Function (PCF), Application Function (AF), Authentication Server Function (AUSF), User Plane Function (UPF), and User Data Management (UDM)) can be virtualized, as it has been proposed for LTE-A, and placed in different virtual machines (VMs) or run as a single bundle in one VM.

Thus, in a scenario where network functions are virtualized, both hardware and software failures assume the same importance, and their reliability shall be guaranteed. Similarly, reliability at service chain level is important to assure proper service availability features to application service platforms deployed by verticals [80, 81].

In this chapter failure detection and recovery mechanism for NG-Core and NG-RAN of the 5G network are presented and experimentally demonstrated.

# 7.1 NG Core Reliability

In 3GPP TR 23.857, different approaches for recovering mobile network functional elements such as: MME failure with/without restart, SGW failure with-/without restart, PGW failure with/without restart and SGSN failure with/without restart. In addition, the following requirements are defined when the above specified services failed:

- The impact or degradation of the service affect on the user shall be minimized;

- The service shall not cause overload or congestion in the Core network;

- The deployment of LTE/EPC and SGs interface shall not impair the delivery of Mobile Terminated CS services (e.g. CS call, SMS).

Thus, when these services are virtualized, there might be a case where the SGW interface may not responsive or VM that deployed the service is not accessible. Hence, it requires to cover the connection by connecting to the hot back or need to restart/initiate the service again.

Furthermore, in the technical specification (TS) 23.007 [82], 3GPP specified different failure detection and recovery mechanisms for core components, including detection of path failure with the help of Echo Request/Echo Response timer messages. Moreover, approaches for recovering failures in a scenario where a mobile network function is virtualised can stem from schemes already proposed for grid and cloud networking. Furthermore, scalable architectures for reliability management are being defined by ETSI NFV [83] and implemented in current open source orchestration frameworks such as Openstack [84].

For example, in [85] resilient schemes for recovering C-RAN failure are proposed based on the concept of access cloud network. In [86] and more deeply explained in Section 7.3, a two-step resiliency scheme orchestrating lightpath transmission adaptation and gNB functional split reconfiguration is proposed. However, the performance of resilience schemes based on the aforementioned approaches once applied to 5G NG-Core have not been fully evaluated so far.

This section shows the capability of recovering vEPC failures by means of a vEPC in "hot backup". Both working vEPC and backup vEPC are deployed in multiple Network Function Virtual Infrastructure Points of Presence (NFVI-PoPs) made available by the federated testbeds belonging to the SoftFIRE project [87].

The goal is to evaluate the Service Recovery Time (SRT), that is the time required to regain user equipment (UE) connectivity, when the proposed resilient scheme is deployed in different NFVI-PoPs.

## 7.1.1 The SoftFIRE Middleware Framework

The SoftFIRE Middleware Framework [88] provides an orchestrated federated virtualisation testbed consisting of component testbeds in multiple countries in Europe. The testbed provides virtualisation platforms controlled by OpenStack [89] and orchestrated by ETSI MANO [83] compliant Open Baton [90] orchestrator. SoftFIRE provides the testbed as an experimentation platform to third parties, which would like to test their 5G applications and virtualisation solutions in a real multi-site testbed. The platform has a Middleware solution based on TOSCA experiment definitions, parsed by an Experiment Manager component. This component invokes relevant sub-managers based on the virtualisation resources requested by an experimenter and provides requested resources (i.e. NFV monitoring, NFV deployment, security as a service, SDN, or reservation of physical devices such as UE or Long Term Evolution-LTE femto-cells). Open Baton then deploys the requested network service defined by the experimenter file, on requested component testbed(s).

## 7.1.2 Considered scenario and the proposed resilience scheme

The considered scenario and the proposed resilience scheme are depicted in Fig. 7.1 and Fig. 7.2 by referring to functional elements of the LTE-Advanced (LTE-A) architecture. The proposed resilience scheme considers a scenario where the vEPC fails (e.g., a virtual machine where the vEPC runs crashes).

Fig. 7.1 shows the two considered vEPC resilience schemes based on vEPCs hot backup deployed in federated NFVIPoPs. The one on the left features two co-located vEPCs (i.e., vEPCa and vEPCb deployed in Surrey 5GIC testbed) while the one on the right features a remote hot backup vEPC (i.e., vEPCr) deployed in a different compute resource available in another testbed (i.e., FOKUS). In the latter case two testbeds will be contemporarily utilized to implement the resilience scheme. In the Surrey 5GIC testbed two different VNFs (i.e., vOAISIM and vEPC) will be implemented by exploiting open source mobile platforms (i.e., OpenAirInterface5g). Here, vOAISIM VNF provides emulation of virtual user equipment

**Figure 7.2:** Proposed scheme experimental evaluation setup

(i.e., vUE) and a virtual gNB (i.e., vgNB) while vEPC will be used to emulate the core network.



**Figure 7.1:** RAN and Core network deployment in SoftFIRE environment

Fig. 7.2 shows the considered scenario and life-cycle event when vEPC VNF fails. Here, when VNFs are deployed, vOAISIM connects with vEPCa, and Zabbix server start monitors the VNFs that are associated corresponding Zabbix agent. Note that each vEPC VNF and vOAISIM VNF deployment contain also Zabbix agent. If the Zabbix server detects an anomaly activity in vEPCa (e.g., overload) or does not receive any status report from vEPCa (i.e., vEPCa crashed) for a predefined period of time (i.e., time to trigger the activity), the Zabbix server check the status of the hot backup vEPC to initiate a recovery procedure. The recovery procedure consist in reconfiguring vOAISIM to connect to the hot backup vEPCr. Upon reconfiguration vOAISIM is able to communicate hot backup vEPCr. Similarly, the experiment also demonstrate the recovery based on the local vEPCb deployed in 5GIC testbed.

## 7.1.3 Experimental Results

The initially considered performance evaluation parameter is the SRT, that is the time required to regain UE connectivity. SRT is measured as the time elapsing between the last ping reply sent by the vEPCa to the vUE before a hot backup remote connection initiation and the detection of the first successive ping reply after successful vUE reconnection with vEPCr (as shown in Fig. 7.3). Ping messages from the vUE to the EPCa are sent every 1ms.

Failure detection is implemented by configuring an action in the Zabbix monitoring server to detect an anomaly activity in vEPCa by monitoring output traffic from gtp0 interface (e.g., interface is not responsive) for a configurable period of time (i.e., 10ms). Once the anomaly is detected, the Zabbix server starts the recovery procedure (as shown in Fig. 7.2) to connect with remote or local hot backup vEPC. Here, the remote hot backup (vEPCr) connection is considered when the active vEPCa failure is detected.

Fig. 7.3 reports the initial vUE is attached to vEPCa both located in Surrey 5GIC testbed, and a Wireshark capture is performed at the vUE interface (172.16.0.2) of the ping messages exchanged by vUE and vEPCa (GTP interface).



**Figure 7.3:** Capture (at the vOAISIM VNF) of ICMP messages exchanged between the vEPCa and vEPCr

When the vEPCa fails (e.g., interface is not reachable) the proposed resilience scheme is activated to connect vEPCr deployed in FOKUS testbed. The timestamp of the Wireshark is measured in seconds. As shown in the Fig. 7.3, once the recovery process is undergoing to connect from vEPCa to vEPCr, only ICMP request messages at the vUE interface are observed. The measured SRT is around 10s. Please note that the evaluated SRT includes also propagation between 5GIC

testbed located at the University of Surrey (UK) and FOKUS testbed located at Fraunhofer (Germany). In addition we are relaunching OAISIMv in order to attach to the hot backup vEPC. However, out latest work showed without service downtime the connection between EPC and gNB can be recovered as described in [91]. The implemented scheme proposed in [91] is able to overcome a single fiber failure in the fronthaul without causing any packet loss, thus ensuring that the DU and CU modules remain connected during the recovery phase.

The proposed demonstration showed the performance of resilience schemes for virtualised mobile network functions. In particular, the demo focused on evaluating the time required to regain UE connectivity when a hot backup vEPC is deployed close to or away from a working vEPC. The demo exploited the remote access to federated testbeds.

## 7.2 NG RAN Reliability

In a several studies, passive optical network (PON) combined with wavelength division multiplexing (DWDM) is the proposed technology to implement the NG-RAN midhaul [92–95]. PON is a cost-effective solution based on a star fiber layout and one passive hub. However, PON does not offer path diversity across the midhaul and a single fiber failure may disconnect one or more DUs from their CU, thus causing service disruption to the UEs connected to the affected DUs. A more reliable fiber infrastructure for the C-RAN makes use of Reconfigurable Optical Add Drop Multiplexer (ROADM) nodes. ROADM nodes are interconnected by fiber spans to form an arbitrary topology.

Any number of fiber failures can be overcome in the midhaul as long as the fiber topology remains connected and the optical circuits carrying baseband signals can be routed (or restored) by the ROADM nodes through working fiber spans. [96,97] provide a more thorough discussion on resilience and reliability in DWDM networks using reconfigurable ROADM nodes. For increased reliability, a DU-CU pair can be connected using two physically disjoint optical circuits, e.g., two circuits that do not make use of the same network fiber span.

Consequently, any single fiber failure can at most disrupt one of the two optical circuits, and the DU-CU pair remains connected through the other-still functioning-optical circuit. The pair of optical circuits can be used by the electronic upper layer to implement protection mechanisms, like 1:1 and 1 + 1. With the former mechanism, data packets are routed over one of the two circuits (pri-

mary), while the other circuit (secondary) is in standby, ready to be used if the primary circuit fails. With the latter mechanism, every data packet is duplicated and concurrently transmitted over the two circuits. One of the two transmitted packets must be removed at the receiver end, while the two circuits are functional.

In the following of this section first a NG-RAN testbed whose midhaul are realized using a highly reliable two layer Ethernet-over-DWDM transport network prototype is described. The programmable optical network (PROnet) testbed (which consists of four ROADM nodes) [98] is configured to provide the midhaul transport functionalities with two built-in fault tolerant mechanisms, i.e., 1:1 + R and 1 + 1 + R, where the 1:1 and 1 + 1 protection mechanisms are implemented at the Ethernet layer and the R (optical circuit restoration) mechanism is implemented at the DWDM layer. The 1:1 protection data plane is implemented using the fast failover table function already available in OpenFlow [99] and Open vSwitch (OVS) [100]. The 1 + 1 protection data plane (which is not readily available in OpenFlow/OVS) was specifically developed by the authors with kernel software and integrated in the PROnet prototype. The protection and restoration mechanisms are software-defined through the PROnet orchestrator.

The second contribution of this section is the application and experimental evaluation of both the 1:1 and 1 + 1 protection mechanisms to achieve reliable midhaul transport networks when using the functional split 7-1 option. As shown in the following sections, both 1:1 and 1 + 1 protection mechanisms suffice to achieve fault-tolerant backhaul transport functionality, which provides uninterrupted service to the UE connected to the NG-RAN even during a fiber outage. However, in the midhaul, the 1:1 protection mechanism's prolonged transport service downtime experienced during the failure-handling procedure is shown to cause disruption of the UE mobile service, resulting in the UE being disconnected from the NG-RAN. In the midhaul, only the 1 + 1 protection mechanism is able to provide a sufficiently fast failure-handling procedure which does not disrupt the UE mobile service from the NG-RAN.

## 7.2.1 PROnet Testbed Configuration

The experimental demonstration is carried out in the PRogrammable Optical network (PROnet) testbed testbed. The PROnet testbed is a 2-layer-Ethernet-over-DWDM-network deployed at and around the campus of the University of Texas at Dallas [101]. Some of the PROnet optical gears are hosted in the OpNeAR Lab,

**Figure 7.4:** PROnet Testbed

next to the Ettus B210 boards and the servers (hosts) used to run the OAI software modules.

A possible configuration of PROnet is shown in Fig.7.4

In this configuration, host's traffic in the network data plane first goes through an Open vSwitch (OVS) controlled with OpenFlow 1.3. The OVS is equipped with multiple 1GE interfaces to connect the various hosts (e.g., servers running the OAI software modules). The data traffic from the hosts is routed by the OVS through a Dell N2048 switch, which is equipped with 1GE and 10GE ports. The Dell switch is used in a static configuration to connect the OVS to the optical muxponder, which offers 10 ports at 10GE, multiplexed together to be transmitted over a single 100Gbps BPSK optical circuit routed at the DWDM layer. The DWDM layer consists of four Cisco NCS 2000 ROADM nodes that are controlled and reconfigured dynamically through NETCONF/ YANG and TL1 interfaces.

The principles of software-defined networking (SDN) are extended in PROnet to control and coordinate the use of equipment at both the Ethernet and the DWDM layers. PROnet is controlled using the PROnet orchestrator which provides two critical functions:

- on-demand data flow provisioning and

- coordinated fault handling at both the Ethernet and DWDM layers.

The PROnet orchestrator is an SDN-enabled multilayer network orchestrator that offers APIs to provision end-to-end network resources with automatic fault-handling features. The PROnet orchestrator discovers the network topology at

both layers, computes two disjoint paths between the endpoints of incoming flow requests, and executes recovery procedures in the event of a network failure, e.g., a fiber cut. By jointly controlling the resources at the two layers, the PROnet orchestrator can efficiently combine restoration and protection mechanisms together, thus achieving a highly reliable multilayer network. Restoration mechanisms are implemented at the optical layer, i.e., when an optical circuit fails, a restoration circuit is automatically established by the PROnet orchestrator. Protection mechanisms are implemented at the Ethernet layer, i.e., two disjoint flows are established ahead of the failure, thus ensuring that at least one flow is readily available and functional in the presence of any single network failure.

Two protection mechanisms and their respective implementations in the PROnet testbed are described next. These are the 1:1 protection mechanism, achieved through the fast failover table protection in OVS, and the 1 + 1 protection mechanism, specifically developed by the authors as a kernel module. Both protection mechanisms can be combined with the optical circuit restoration mechanism in PROnet, thus achieving highly reliable network services like 1:1 + R and 1 + 1 + R, which can overcome subsequent network element failures.

The connection between the NG Core and the CU, i.e. the backhaul, is implemented through direct Ethernet cable. The connection between the DU and the CU, i.e., the midhaul, is implemented through optical circuits that are established in the PROnet testbed. The optical circuits are arbitrarily interrupted during the experiment to evaluate the effects of such interruption on the midhaul services.

## 7.2.2  1:1 + R protection and restoration mechanism

Fig. 7.5 shows an example that illustrates how fault handling is implemented using the fast failover table protection. Fault handling here is a combination of a proactive approach (protection) and a reactive approach (restoration). When an Ethernet flow between two hosts is requested, the PROnet orchestrator allocates protection resources, i.e., two separate and physically disjoint network paths are allocated. In the figure, the primary path is colored in green (carrying the flow) and the secondary path in red (in standby and not used to transmit the flow). In the presence of a failure in the primary path, the flow is switched from the primary path to the secondary path at the Ethernet layer. This flow switching is achieved using the OpenFlow feature called *fast failover table* [102]. This feature is triggered by a local detection of the failure, i.e., the OpenFlow-enabled Ethernet

**Figure 7.5:** 1:1 + R protection and restoration mechanism

switch detects the link failure locally and promptly forwards the disrupted flows to the other (secondary) link.

In particular, the switch OVS A, initially forwards the traffic using the first entry in the fast failure table (located at the OvS A). Therefore the switch OVS A uses the port 2 (i.e., connected to primary path) until its status is changed to down, i.e. when the primary path (i.e., green link) is not operational. At this occurrence, the flow is temporarily forwarded using the second entry in the fast failure table and therefore using the OVS A port 3 (i.e., connected to the secondary link). Please note that the failure is manually created by putting the OVS A interface 2 down. The flow then, remains operational despite the fault that disrupted the primary path in the network. Regarding the DWDM layer, please note that ROADM nodes are transparent to the failure being them configured with two ports: a first ROADM nodes port is connected to the primary path (i.e., green path) while a different ROADM nodes port is connected to the secondary path (i.e., red path). Therefore, based on the port that is receiving traffic, the ROADM node W3 forwards the received traffic to the corresponding output port.

The bottom part of the figure shows an example of fast failover table entries. At the Ethernet switch, the incoming flow is matched against source and destination MAC address (IP address and port number can be used as well) to a group table entry. Each group table entry has two forward options represented by the Bucket IDs.

Note that, however, the flow now relies on the secondary path only, which is not designed to overcome a second subsequent fiber outage. To address this issue,

the PROnet orchestrator automatically performs a restoration procedure at the DWDM layer immediately after the fault detection. The objective of the restoration procedure is to compute and establish a restoration optical circuit to rebuild the disrupted primary path. A critical routing requirement for the restoration circuit is to be physically disjointed from the secondary path. This condition guarantees that any subsequent fiber span failure would affect at most only one of two paths: either the newly established restored path or the secondary path. When the optical circuit restoration procedure is complete, the primary link interface status at the switch is changed to up and the flow is forwarded to the newly established restored path. Notice that the Ethernet switch is not aware of the change of primary path that took place at the DWDM layer, and the flow is now protected (again) against any single and subsequent fiber span outage.

It is not uncommon to experience some packet loss during the protection procedure, as the switch takes time to detect the link outage and forward the flow over the secondary path.

### 7.2.3  1 + 1 + R protection and restoration mechanism

An Ethernet 1 + 1 protection mechanism was recently added to the PROnet testbed. This mechanism is designed to overcome a single fiber failure in the network without causing any packet loss.



**Figure 7.6:** 1:1 + R protection and restoration mechanism

Fig. 7.6 shows an example that illustrates how the 1 + 1 + R mixed protection and restoration mechanism work. Two disjoint paths are computed and provisioned by the PROnet orchestrator. Every packet in the flow is duplicated at the

ingress switch and the two copies are continuously forwarded, one to the primary and one to the secondary path. At the other end of the two paths, the egress switch receiving the duplicated flows must guarantee that only one copy of each packet is delivered to the host. A Linux kernel software module had to be custom-developed to implement the 1 + 1 protection mechanism. The newly developed kernel software operates as briefly described next.

The ingress flow, i.e., the flow that is entering the PROnet optical network, is copied to two output interfaces by the ingress OVS switch, using the *all group table function* [102]. At the egress OVS, the two copies of the same flow arriving from the two PROnet optical circuits are processed in real time in order to remove packet duplicates before forwarding to the receiving end. The egress OVS processes the packets from the two incoming flows using netfilter hooks. More specifically, the hook used is the pre-routing hook. Each time a packet is captured by the hook, the following actions are performed:

- a unique signature is created for the packet using some of the packet field(s);

- if the packet signature is found in the checklist (i.e., a copy of this packet was received earlier), the packet is discarded and the signature value is removed from the list. Indeed, the module assumes that a packet can only be duplicated and no packet will be received more than twice.

- else (this is the first copy of the packet being received) the packet is passed onto the switch which performs the necessary forwarding and the packet signature is stored in the checklist.

The software also removes the oldest signatures from the checklist when the list exceeds a pre-defined number of entries, N. This additional feature is required to prevent the checklist from growing in size unbounded under certain conditions. For example, when one of the two copies is not received by the module (e.g., one of the optical circuits is disrupted), the signature of the packet that is received only once would remain in the checklist for ever. The value for N is programmable and can be computed to account for the transmission rate, packet size, and difference in signal propagation time between the two disjoint paths. The results discussed in the next section are obtained setting N = 100 , which proved to be sufficient large to ensure correct operation.

Note that the packet signature must be computed to be unique for each received packet and its copy. For example, when transmitting packets over the midhaul

using functional split option 7-1, the signature is computed using data in the RLC packet, i.e., the signature is the frame number, subframe number, and symbol number fields, as these three parameters are combined to form a unique triplet for each packet that is transmitted on the midhaul. As such, the calculation of the signature does not contribute any significant time to the midhaul latency, i.e., of signatures and comparison with signatures stored in the checklist takes less than 1 $\mu$s.

The 1 + 1 mechanism can be combined with the R (restoration) mechanism to achieve the 1 + 1 + R protection and restoration mechanism by applying a technique similar to the one already described for the 1:1 + R mechanism.

## 7.2.4 Experimental Results

The testbed described in the previous sections was used to conduct a number of experiments. These experiments were designed to assess the reliability of the midhaul networks. Reliability is defined here as the ability of the testbed to maintain the UE connected to the NG-RAN when a network outage, e.g., a fiber span in PROnet is disconnected, is artificially introduced.

The midhaul reliability is first tested using the 1:1 fast failover table protection mechanism. Network resources are first provisioned through the PROnet orchestrator, e.g., primary and secondary optical circuits, along with OVS fast failover tables. Then the C-RAN modules are manually started. Once the NG-RAN is fully operational and the UE is connected, a fault is introduced by manually disconnecting one of fiber spans used by the primary circuit. As soon as the network outage is introduced, we observed that the NG-RAN stops working and the UE is disconnected from the NG-RAN despite the fact that the fast failover table protection recovers from the fiber span failure within a few tens of milliseconds. This is due to the strict network latency and jitter requirements that must be fulfilled between the CU and DU. As reported in [103–105], when CPRI is utilized to transport lower-layer split options (e.g., option 8 and option 7-1), it needs to support delay within 100 $\mu$s and jitter within 65 ns, because of the time-sensitive I/Q data that is being transmitted. An additional drawback of the fast failover table protection mechanism is the fact that this protection mechanism does not guarantee delivery of all the transmitted packets over the midhaul, i.e., at the moment of the fiber outage occurrence and until the outage is circumvented by the failover table, the transmitted packets are dropped at the switches attached to the primary cir-

**Figure 7.7:** Wireshark capture (CU side) capturing echo requests and replies between the CU and the DU

cuit affected by the failure and are not delivered to the intended destination. The current C-RAN implementation in OAI is sensitive to any packet loss experienced on the link connecting the CU and DU and does not appear to be able to recover from said packet loss. Manual restart of the DU and CU modules is therefore necessary at this point to reconnect the UE to the NG-RAN.

A second experiment is conducted, this time by applying the newly implemented 1 + 1 protection mechanism. Network resources are first provisioned through the PROnet orchestrator. Then the NG-RAN modules are manually started. Once the NG-RAN is fully operational and the UE is connected, a network outage is introduced by manually disconnecting one of the fiber spans used by either the primary or secondary optical circuit. When operating with the 1 + 1 protection mechanism in the PROnet testbed, no loss of data packets is detected even, while the fiber span is disconnected and the 1 + 1 protection mechanism must account for the loss of one of the optical circuits. Fig. 7.7 shows the Wireshark output, while monitoring the CU network activities. The CU server (IP address 192.168.0.134) is set to transmit periodic echo request packets to the DU server (IP address 192.168.0.136), at intervals of about 200 ms. The row highlighted in orange indicates the first echo request (sequence number 404) that is sent after the occurrence of the outage. For determining the exact moment in which the 1 + 1 protection mechanism is called into action, we applied a pragmatic procedure to this experiment as follows. The kernel module that filters packets at the egress node dropping duplicate packets is applied only to NG-RAN data packets between the DU and CU, while ICMP, i.e., ping packets are not removed at all. It should be noted that when using IF4.5

in OAI, the interface between DU and CU transfers OFDM symbols with frame, subframe, and symbol count for both the uplink and downlink, as well as PRACH packets. As a result of this arbitrary procedure, when pinging the server hosting the CU from the server hosting the DU, both echo requests and echo reply packets are duplicated by the 1 + 1 protection mechanism, resulting in 4 echo reply packets being detected with the same sequence number. Once the outage disrupts one of the two optical circuits, both request and reply packets are received only once, as the echo packets cannot make it through the failed (bidirectional) optical circuit. By inspecting the Wireshark output, it is straightforward to identify the exact moment the optical circuit is disrupted. The contiguous sequence numbers reported in the trace reveal that no packet is lost during the experiment, including the moment when the optical circuit is disrupted. The UE remains connected to the NG-RAN thanks to the lossless 1 + 1 protection mechanism.

# 7.3 Lightpath Recovery and Functional Split Adaptation in NG RAN

The establishment of 5G technology will have also a deep impact on the design, control, and management of metro networks [106–108]. The dimensioning of metrorings will be driven by the maximum latency admitted by 5G services and verticals. The network capacity will increase to support several high-bandwidth midhaulbackhaul connections. Thus, the employment of transponders based on high-spectral efficient multilevel modulation formats (e.g., polarization multiplexing 16 quadrature amplitude modulation [PM-16QAM]) should penetrate this market. Moreover, proper programmability and control of transmission systems and switching will enable dynamic provisioning of connectivity and its survivability.

A way considered by operators and vendors to reduce the costs of an optical network is the reduction of margins [109]. Indeed, margins to the quality of a transmission are typically adopted to consider uncertainties of physical layer models and device aging. This is a pessimistic but conservative approach that enables uninterrupted service of connections during the whole life of the network but that could also result in overestimation of the number of regenerators in regard to costs. Reducing margins can decrease costs (e.g., of regenerators) but can also increase the probability of experiencing soft failures [110] (i.e., degradations of the connectivity resulting in a bit error rate [BER] increase over the acceptable thresholds)

due to model uncertainties and aging. Differently from hard failures (e.g., link cut), where only rerouting enables traffic recovery, soft failures can also be overcome by adopting a more robust transmission along the degraded path, e.g., by adapting the modulation format (e.g., changing from PM-16QAM to the more robust polarization multiplexing quadrature phase shift keying [PM-QPSK]) or by increasing code redundancy. Hereafter, "lightpath adaptation" refers to operations involving the change of transmission settings of a lightpath, as the adaptation of the modulation format, of the rate, or of the code redundancy. Lightpath adaptation can be particularly fast and, in some cases, hit-less [111] (e.g., change of rate).

However, modulation format change or a code redundancy increase at a fixed baud rate implies an information rate reduction. Thus, if recovery is performed through a modulation format or code adaptation, part of the traffic can be promptly recovered along the same path, while the other has to be rerouted or, alternatively, suppressed if the service class admits a bit-rate reduction [112].

Therefore, it is expected that a metro network will carry a converged mix of traffic, including midhaul, backhaul, and non mobile traffic, requiring large capacity. In this scenario, guaranteeing reliability is of paramount importance. The failure scenarios can be several (i.e., virtual machine/container failure, optical link failure, etc.), but this section focuses on the failure of the midhaul connection carrying the virtual link between a virtualized DU (vDU) and a virtualized CU (vCU).

In particular, a lightpath, which carries the midhaul connection, soft failure (i.e., a degradation of the quality of transmission) is considered. In this scenario, upon a failure, the connectivity between vDUs and vCUs can be easily recovered by changing the modulation format or increasing the code redundancy. However, the consequent rate/capacity reduction might not be acceptable by the midhaul, implying the rerouting of the original lightpath along a different route. Nevertheless, such rerouting might fail if the available capacity on the other direction of the ring is not sufficient. This section presents and evaluates a two-step recovery scheme, stemming from the one presented in [113], orchestrating lightpath adaptation and gNB functional split reconfiguration to recover the vDU-vCU connectivity while fulfilling the midhaul capacity requirements. Although resilient schemes for recovering failures when network functions are virtualized have been already investigated [114] or can be based on previous research on cloud [115] and grid computing [116], the originality of the proposed scheme consists in exploiting functional split flexibility to improve vRAN reliability. The scheme first recovers the soft failure by lightpath adaptation or lightpath rerouting, and, if the provided

| Split Option | Required Bandwidth | Maximum Allowed One Way Latency |
|:---:|:---:|:---:|
| **Option 2** | DL: 4 Gb/s; UL: 3 Gb/s | 10 ms |
| **Option 7a** | DL:10.1-22.2 Gb/s UL:16.6-21.6 Gb/s | 250 $\mu$s |
| **Option 7b** | DL:37.8-86.1 Gb/s; UL:53.8-86.1 Gb/s | 250 $\mu$s |
| **Option 7c** | DL:10.1-22.2 Gb/s; UL:53.8-86.1 Gb/s | 250 $\mu$s |
| **Option 8** | DL: 157.3 Gb/s; UL: 157.3 Gb/s | 250 $\mu$s |

**Table 7.1:** Midhaul Requirements

capacity is not sufficient, it modifies the functional split between vDUs and vCUs, so that it can be supported by the new lightpath capacity.

The results, collected via simulation and experimental evaluation, show that, at the beginning, the vDU-vCU virtual link incurs a graceful degradation due to lightpath adaptation. Then, upon functional split reconfiguration, the vDU-vCU virtual link is recovered to fully support the required capacity. The overall recovery time is in the order of a few seconds if several vDU and vCU functional splits are pre-deployed.

## 7.3.1 Considered Scenario and Two-Step Recovery Scheme

This section summarizes the considered scenario and the proposed two-step recovery scheme. The city of Milan is considered a sample metropolitan city. The city metropolitan area covers about 200 square kilometers. By assuming an antenna density of one antenna per 2 square kilometers, 100 antenna sites provide connectivity to the entire area. If split Option 2 is utilized as a functional split, the capacity required by the connections between DU and CU is about 4 Gb/s in each direction, as reported in Table 7.1, from 3GPP TR 38.801. This connection can be supported by a 10 GbE link between DU and CU. By multiplying each antenna midhaul requirement by the number of antenna sites, the overall required capacity sums up to 1 Tb/s, as summarized in Tab. 7.2.

Based on the data reported above and on practical deployment considerations [117], the architecture depicted in Fig. 7.8 is considered. Antenna sites are connected to optical switches that form an optical metro-ring network. The network

| Milan, Italy | |
|---|---|
| **Metropolitan area** | 200 $km^2$ |
| **Number of antenna sites** | 100 |
| **Average antenna density** | 0.5 $km^{-2}$ |
| **Functional split Option** | 2 |
| **Required capacity per antenna site** | 4 Gb/s per direction (UL/DL) |
| **Overall midhaul required capacity** | 1 Tb/s |

**Table 7.2:** Metropolitan Area Capacity Requirements

architecture connecting the antenna site and metro-ring switch can be varied; two possible solutions are point-to-point connections or a next-generation passive optical network (NG-PON2). DUs and CUs are virtualized. vDUs are placed either at the antenna sites or in their vicinity (connected to the same metro switch). vCUs are placed in another data center connected to one of the metro-ring switches. The connectivity between vCU and vDU is implemented by means of lightpaths routed in the optical metro-ring network with an 8 km radius. Such radius implies a circumference of about 50 km with a maximum latency of about 250 $\mu$s (considering propagation delay only).



**Figure 7.8:** Metro midhaul architecture

The proposed scheme (shown in Fig. 7.8) is in turn, derived from the ETSI NFV-MANO architecture [118]. The network function virtualization orchestrator (NFVO) is responsible for orchestrating vDU, vCU, and network resources and for selecting the proper functional split. The virtual infrastructure manager (VIM) and virtual network function manager (VNFM) are responsible for provisioning the required compute, storage, and network resources and deploying vDU and the vCU according to the functional split selected by NFVO. The SDN controller is responsible for the control of vDU-vCU connectivity, i.e., for the configuration of the optical metro segment. The SDN controller is enhanced with a network monitor, monitoring the lightpath status.

Fig.7.9 shows a flow chart of the proposed two-step recovery scheme. Upon detection of a soft failure, lightpath adaptation is triggered.

**Figure 7.9:** Two-step recovery scheme flow chart

Lightpath adaptation is performed by the SDN controller, which configures the transmitter and receiver with the new configuration settings. Lightpath adaptation can be based on a modulation format adaptation, which can be performed by a digital-to-analog converter (DAC), typically employed in transponders based on coherent transmission, and is now expected to also penetrate the metro network market [119]. Upon adaptation, if the functional split rate is still supported by the lightpath, the original functional split is recovered and maintained. Otherwise, another path is searched to recover the original rate. If no path is found or, generally, the functional split cannot be supported, the functional split reconfiguration

is triggered.

A functional split is sought starting from the highest split option (i.e., lowest functional layer split) that could be carried by the adapted lightpath. If successful, the vDU-vCU connectivity is recovered. Otherwise, it is lost. The decision to start the search from the highest split option is to minimize the degradation on the wireless transmission performance.

As an example, the working lightpath (red solid line in Fig. 7.8) connecting vDU and vCU is monitored by the network monitor. Such a monitor reveals or even anticipates degradation in the quality of transmission of the working lightpath, and it notifies the SDN controller, which triggers the modification of the lightpath modulation format or the code redundancy increase (i.e., green dashed line) if needed. If the resulting lightpath rate is not capable of carrying the original functional split (i.e., vDU1s1 and vCU1s1), the SDN controller attempts to establish another lightpath along the opposite direction of the ring. If unsuccessful, the NFVO is notified, and it triggers the modification of the functional split to one requesting a lower data rate (i.e., vDU1s2 and vCU1s2), and the vDU-vCU connectivity is recovered along the original path with an adapted lightpath.

## 7.3.2 Simulation Scenarios

Simulations are carried out to evaluate the amount of midhaul interfaces that can be recovered through transmission parameter adaptation (thus, without rerouting) in a metro network. A custom-built event-driven C++ simulator is utilized, and a five-node ring topology is considered. Each link is 10 km long, not needing in-line amplifiers. Optical amplifiers are assumed as boosters. In the metro network, connection requests, aggregating 5G services, are modeled as a Poisson process (e.g., emulating a small cell on-off process). The holding time of each connection is exponentially distributed with average $1/\mu = 5000$ s. The traffic load offered to the network is expressed as $\lambda/\mu$, where $1/\lambda$ is the mean inter-arrival time of connection requests. Transponders supporting 200 Gb/s PM-16QAM and 100 Gb/s PM-QPSK are considered. The bit error rate (BER) of both PM-16QAM and PM-QPSK is computed through the optical signal-to-noise ratio (OSNR) assuming negligible nonlinear effects given the limited distances of the metro network [120–122]. A BER lower than $10^{-3}$ is assumed as acceptable. In particular, 200 Gb/s PM-16QAM is considered acceptable with an OSNR higher than 20.5 dB, while the model in [122] is adopted for 100 Gb/s PM-QPSK. Soft failures are ran-

domly generated on a single link. We assume that a 200 Gb/s lightpath carries an aggregation of Option 8 functional split traffic, while a 100 Gb/s lightpath carries an aggregation of Option 7 functional split traffic. It is assumed that Option 8 and Option 7 midhaul interfaces require 253 Mb/s and 160 Mb/s, respectively; such values were obtained experimentally by using OAI as in [123]. Thus, a 200 Gb/s lightpath carries around 790 Option 8 midhaul interfaces, while a 100 Gb/s carries around 625 Option 7 midhaul interfaces. The proposed two-step recovery scheme is compared with rerouting. When rerouting is applied, if a lightpath is impacted by the soft failure (BER above the threshold), rerouting is performed while keeping unchanged 200 Gb/s (and modulation format), thus unchanged Option 8. Two-step recovery and rerouting are compared in terms of supported midhaul interfaces after the soft failure. Results are recorded until the confidence interval of 5% at 95% confidence level is achieved.

### 7.3.3 Experimental Scenarios

The architecture depicted in Fig. 7.8 is implemented in the a simplified version shown in Fig. 7.10.



**Figure 7.10:** Scenario considered for the experimental evaluation

Here, the evolved packet core (EPC) is deployed in PC 1, and CUs and DUs are virtualized in PC 5 and PC 6, respectively (see Tab. 4.1 in Chapter 4 for more details on the considered machines). As described in 5, openair-cn is utilized as EPC while the functional splits implemented by the openairinterface-5g platform. The functional splits considered in this section are the Option 8 and Option 7a.

**Figure 7.11:** Lightpath transmission adaptation FSM

The midhaul link between vCU and vDU is assumed to be realized by means of an optical metro-network, including sliceable transponders capable of adapting the modulation format and the spectrum occupation according to the required bit-rate and path length. The transmitters (TXs) support different modulation formats (i.e., QPSK, 8QAM, 16QAM) and different baud rates (i.e., 28 and 32 Gbaud in relation to the forward error correction [FEC code rate adopted]). The optical data plane is emulated considering BER values collected through measurements in the setup described by [124]. Each node is equipped with a network configuration protocol (NETCONF) agent, developed using the ConfD framework, in order to enable the SDN paradigm. The SDN controller relies on NETCONF protocol to perform the configuration and the monitoring of the optical nodes. By means of NETCONF edit-config, the SDN controller configures the traversed ROADMS and the considered transponders.

When a lightpath is configured, the receiver periodically reports monitored parameters, including pre-FEC, BER, and OSNR. All the receiver transponders are configured with specific finite state machines (FSMs), in order to automatically adapt the lightpath transmission rate in the case of known events. The considered FSM is shown in Fig. 7.11. As an example, if the lightpath operates at "State 1", but the BER exceeds a threshold, a modulation format adaptation is triggered, and the transition to "State 2" is performed. Conversely, if the lightpath operates at "State 2" but the BER is below a threshold, a modulation format adaptation is triggered, and the transition to "State 1" is performed.

In our experiments, to perform modulation format adaptation, we exploited the scheme proposed in [125, 126] based on FSM for enabling fast reconfiguration. In this way, if the value of BER monitored at the receiver exceeds the configured

threshold, the automatic reconfiguration of TX and RX can be performed locally without involving the SDN controller. The portion of the FSM (State 1) to be adopted for the automatic reconfiguration in the case of monitored BER greater than 0.002 is the following.

**Listing 7.1:** Finite state machine configuration

```
1
2  <state>
3  <id>1</id>
4  <description>State1</description>
5  <events xmlns:ev="http://sssup.it/events">
6    <event>
7      <name>BER</name>
8      <type>ev:ON_CHANGE</type>
9      <threshold−param>0.002</threshold−param>
10     <check−operator>GT</check−operator>
11     <reaction>
12      <operation>
13        <id>1</id>
14        <type>SIMPLE_OP</type>
15        <simple>
16         <local−address>10.1.1.3</local−address>
17         <remote−address>10.1.1.2</remote−address>
18         <execute−local>
19            <transponder xmlns="http://sssup.it/transponder"
20              xmlns:nc="urn:ietf:params:xml:ns:netconf:base:1.0">
21             <subcarrier−module>
22               <subcarrier−id>1</subcarrier−id>
23                <config>
24                  <bit−rate>112</bit−rate>
25                  <baud−rate>28</baud−rate>
26                  <modulation xmlns:mf="http://sssup.it/modulation−formats">mf:pm−qpsk
                        </modulation>
27                </config>
28             </subcarrier−module>
29            </transponder>
30         </execute−local>
31         <execute−remote>
32            <transponder xmlns="http://sssup.it/transponder"
33              xmlns:nc="urn:ietf:params:xml:ns:netconf:base:1.0">
34             <subcarrier−module>
35               <subcarrier−id>1</subcarrier−id>
36                 <config>
```

```
37            <bit−rate>112</bit−rate>
38            <baud−rate>28</baud−rate>
39            <modulation xmlns:mf="http://sssup.it/modulation−formats">mf:pm−
                 qpsk</modulation>
40          </config>
41         </subcarrier−module>
42        </transponder>
43       </execute−remote>
44       <next−state>2</next−state>
45      </simple>
46     </operation>
47    </reaction>
48   </event>
49   </events>
50 </state>
```

A similar configuration file is utilized for the transition to State 2. The reconfiguration consists of adapting the lightpath with a baud rate of 28 Gbauds and a more robust modulation format (i.e., PM-QPSK) at the TX (remote agent) and the RX (local agent). The FSM presents another state (i.e., State 2), designed for the case when the value of BER undergoes a threshold of 0.00004. In fact, in this occurrence, the lightpath is reconfigured with a baud rate of 32 Gbauds and a modulation format with higher bit rate (i.e., PM-16QAM). Moreover, the NETCONF agent acting in the receiver has been extended with a specific NETCONF notification stream, which is able to provide all the required information to the subscribed clients. In our case, the VNFM/VIM performs the subscription to this stream, to receive an update once the lightpath is adapted. Indeed, VIM manages the network function virtualization infrastructure (NFVI) while VNFM manages the virtual network functions. To further speed up recovery operation, in this study, several vDU and vCU types featuring different functional splits are already installed (i.e., containers with the respective functions are already allocated), and they are activated upon request. The number of vDUs and vCUs that could be allocated depends on the compute and storage resource of the data centers. As shown in [127], Docker container-based virtualization allows a higher midhaul latency and jitter budget than virtualization based on virtual machine. Thus, Docker container-based virtualization is considered in this experimental evaluation.

As shown in Fig. 3, PC 5 contains a Docker with two containers to deploy the two vCUs (i.e., vCU1 and vCU2) with two different functional split options. Here, the vCU1 and vCU2 are used to build and run for Option 8 and Option 7a, respec-

tively. Similarly, the vDU functions with two different functional split options are hosted in PC 6 (i.e., vDU1 and vDU2) to build and run Option 8 and Option 7a, respectively. The USRP B210 is attached to PC 6. The Huawei E3372 dongle is utilized as UE. The dongle is connected to PC 7 and connected to the USRP B210 device through SMA cables with 20 dB of attenuation in the middle of the link.

The VIM/VNFM is emulated by a shell script that is triggered by the NFVO when functional split change is required. The main workflow is highlighted in Fig. 3. At the beginning, the VNFM/VIM and the SDN controller perform the setup of the network service, activating vCU1 in PC 5, vDU1 in PC 6 and deploying the lightpath interconnecting them (three nodes traversed with modulation 16QAM and baud rate 224 Gbauds) with Option 8. In Step 1, the receiver (RX) detects a value of BER ¿ 0.002, and consequently, the FSM is applied: TX (Step 2) and RX (Step 3) are reconfigured according to the FSM configuration. Once the reconfiguration is complete, a NETCONF notification is sent (Step 4), informing the VNFM/VIM of the lightpath adaptation. Then, the VNFM/VIM performs the activation of vCU2 and vDU2 with Option 7a (Step 5).

The considered performance evaluation parameter is the gNB functional split reconfiguration time (FSRT), here defined as the time elapsing between the last ping reply sent by the EPC to the UE and the detection of the first successive ping reply with the requested functional split option at the UE. The FSRT measurement is performed as follows:

1. Initially, the vRAN setup is set to run on vDU1 and vCU1 with functional split Option 8;

2. the ping is continuously run between the UE and the EPC with packet interval of 1 ms;

3. the reconfiguration of the functional split request command is sent by the VNFM/VIM to vDU2 and vCU2 to initiate the requested functional split (i.e., Option 7a).

## 7.3.4 Results

In this section are firstly presented the simulation results and the experimental ones.

**Figure 7.12:** Mean overall recovered rate through transmission parameter adaptation versus offered traffic load, with soft failures introducing an OSNR penalty of 2 dB.

Fig. 7.12 shows the number of supported midhaul interfaces before and after the soft failure versus the offered traffic load, assuming that each soft failure introduces an OSNR penalty of 2 dB. While traffic load increases, with rerouting, the lack of available network resources on alternate paths causes the failure of lightpath rerouting; thus, several midhaul interfaces are no longer supported.Conversely, the proposed two-step scheme exploits lightpath adaptation along the same route, and the change of modulation format makes the transmission feasible even with the OSNR degradation; thus, recovery is not blocked for lack of network resources and two-step recovery outperforms rerouting. The plot also shows the number of supported Option 8 and Option 7 interfaces after the failure in the case of two-step recovery. The former belongs to the lightpaths that are not impacted by the soft failure. At high loads, rerouting approaches the Option 8 interfaces with two steps because all lightpath reroutings are blocked for lack of resources; thus, only Option 8 traffic not impacted by the failure is present in the network. The number of supported Option 7 interfaces after the failure slightly increases with load because more lightpaths are impacted by the soft failure; thus, the change of functional split is more likely exploited.

Fig. 7.13 shows the number of supported midhaul interfaces before and after the soft failure versus the OSNR penalty with a traffic load of 500 Erlang. The to-

tal number of supported interfaces after failure decreases with the OSNR penalty. Indeed, the number of lightpaths impacted by the soft failure increases with the OSNR penalty because the larger the penalty the higher the probability of passing the BER threshold of $10^{-3}$. Thus, in the case of rerouting, more lightpaths contend available network resources on alternate paths, while, in the case of two-step recovery, more lightpaths must rely on bit-rate reduction.



**Figure 7.13:** Mean overall recovered rate through transmission parameter adaptation versus soft-failure OSNR penalty, with a traffic load of 100 Erlang.

Regarding the experimental results, Tab. 7.3 shows the average values of the overall transport network reconfiguration time in the considered scenario for lightpath adaption. The value of TX *reconf* reports the time required to reconfigure the TX. The value of the RX *reconf* indicates the time required for both TX and RX reconfiguration. Then, the value at the NETCONF *notify* is the overall time required for the reconfiguration, including the NETCONF notification to the VNFM/VIM related to the lightpath adaptation. Two different lightpath adaption scenarios are considered such as:

1. from high bit rate to low bit rate (i.e., from 16QAM/32Gbauds to QPSK/ 28 Gbauds);

2. from low bit rate to high bit rate (i.e., from QPSK/28 Gbauds to 16QAM/32 Gbauds).

| Lightpath Adaptation | TX Reconf Time [ms] | RX reconf Time [ms] | NETCONF Notify Time [ms] |
|---|---|---|---|
| **From high bit rate to low bit rate** | 145.83 | 250.39 | 294.51ms |
| **From low bit rate to high bit rate** | 166.75 | 272.83 | 315.64 |

**Table 7.3:** Transport Network Reconfiguration Time

These two different scenarios are based on the pre-defined parameters, as shown in Fig. 7.10. The overall transport network reconfiguration from high bit rate to low bit rate is around 295 ms, while from the low bit rate to high bit rate is around 316 ms, as shown in Tab. 7.3.



**Figure 7.14:** Capture, at the UE, of ICMP messages exchanged between the UE and the EPC

Upon receiving the NETCONF notification about the lightpath adaptation as shown in Fig. 7.10, the VNFM/VIM performs the activation of vCU2 and vDU2 with Option 7a. Fig. 7.14 illustrates the wireshark capture of ping messages at the UE during the gNB functional split reconfiguration. The ping messages (i.e., ICMP request and reply messages) are exchanged between the UE (dongle), whose IP address is 192.168.8.1, and the EPC, whose GTP interface IP address is 172.16.0.1. Notice that the vRAN setup is established and running with functional split Option 8, initially. The timestamp of the wireshark is measured in seconds. As shown in Fig. 7.14, the last ICMP reply message is received by the UE at timestamp 197.930822 s before functional split Option 7a is triggered. Upon functional split reconfiguration being triggered, only ICMP request messages at the UE can be observed. The first successive ICMP reply message from the EPC to the UE is received at timestamp 210.392784 s, showing the successful reconfiguration of the functional split Option 7a. Thus, the time elapsing between the last ICMP

reply message and the first successive ICMP reply message at the UE (i.e., the FSRT) is about 12 s. This obtained FSRT values also includes a 2 s sleep time to synchronize the midhaul interface between OAI DU and OAI CU during Option 7a configuration, and the shell-based VIM controller contribution of about 1 s to enter into the containers to run the requested functional split option.

## 7.4 Conclusion

This Chapter described a testbed conceived to test reliability mechanisms applied to both backhaul and midhaul in NG RAN.

In Section 7.1, a resilience scheme for the NG Core is proposed. It is shown the performance of resilience schemes for virtualised mobile network functions. In particular, the Section 7.1 focused on evaluating the time required to regain UE connectivity when a hot backup vEPC is deployed close to or away from a working vEPC. The study exploited the remote access to federated testbeds.

Section 7.2 described a testbed conceived to test reliability mechanisms applied to midhaul in NG RAN. The use case in this study was that the NG RAN functional split option 7-1 was provided by OAI running on a number of Linux servers. The midhaul were implemented using either a simple Ethernet cable or the two-layer Ethernet-over-DWDM PROnet. Two reliability mechanisms were discussed and applied to the midhaul to improve its reliability, defined as the ability of the network not to disconnect the UE from the NG RAN in the presence of fiber span outage. The first mechanism makes use of 1:1 protection at the Ethernet layer combined with an optical circuit restoration at the DWDM layer. The second mechanism makes use of 1 + 1 protection at the Ethernet layer combined with an optical circuit restoration at the DWDM layer. These two combined protection and restoration mechanisms (1:1 + R and 1 + 1 + R) are coordinated by the software-defined network (SDN) PROnet orchestrator to overcome any single fiber span outage, which may be followed by other subsequent outages. More specifically, the implementation of the 1 + 1 protection mechanism required the development of a new kernel software module, which runs in Open vSwitch (OVS) and achieves zero packet loss even during the occurrence of a fiber outage. Experimentally, it has been determined that in the midhaul, the 1:1 protection mechanism was failed to keep the UE connected to the C-RAN, due to the prolonged traffic disruption and the resulting burst of lost data packets. However, by applying the 1 + 1 protection mechanism to the midhaul, the UE remained connected to the C-RAN even after

the fiber outage occurrence.

Finally, Section 7.3 proposed a two-step scheme for recovering virtualized distributed unit (vDU) and virtualized central unit (vCU) connectivity upon midhaul failure. In particular, a single soft failure of a lightpath interconnecting vCU and vDU is considered. The main novelty of the proposed scheme consists of complementing lightpath adaptation with the possibility of flexibly changing the vCU and vDU functional split. This approach allows the increase of the number of the recovered connections when spare network resources are scarce. Simulation results showed that the proposed scheme allows us to recover almost all the midhaul vDU-vCU connections along the same path where working lightpaths were routed, and it largely overcomes the performance of lightpath rerouting, especially when the network load is high. Experimental results showed that the recovery time experienced when functional split change is triggered is in the order of tens of seconds.

# 8 Slicing in 5G networks

In traditional network infrastructure, network functions (NFs) are deployed as physical proprietary devices (software and hardware). In fact, to deploy a NF, a specialized physical appliance was needed with pre-installed closed software. With the evolution of the network, due to the introduction of SDN and the concept of NFV, it is possible to deploy NFs as virtualized elements for a flexible and efficient utilization of the infrastructure.

This vision has brought ETSI in defining the specification of a framework in support of NFV management and orchestration. The framework objective is to support VFN operation across hypervisors and computing resources, that means the spawning of virtual function in one or more shared Data Centers (DCs). Moreover, it specifies how to orchestrate and manage the life-cycle of physical and virtual resources. According to [128], the framework does not propose any specific implementation and it is described at a functional level.

One of the main challenges in the field of networking is to meet the user requirements. With the 5G technology, we consider as users the verticals. A vertical is a socio-economic subject that offers services: from the e-health to the e-entertainment, from the automotive to the pure telecommunication. Each service offered has stringent requirements in term of network reliability, latency and transmission bandwidth. To accommodate such services it is needed to deploy a dedicated portion of network (slice) and instantiate the components of the services in distributed clouds or to the RAN edge (i.e. to the MEC).

This operation are performed by an architectural element that has the role of manage and orchestrate the network services. The reference architecture for this element has been proposed by the ETSI with the MANO framework. The framework can be extended according to specialized functional requirements.

Its main objective is to build up adaptive customized services. In fact, each service responds to specific characteristics and, since the resource utilization (both network and cloud) is variable, it is necessary to constantly adapt them.

In this context, the H2020 5G-TRANSFORMER project (5GT) [22] aims at trans-

forming today's rigid mobile transport networks into a flexible SDN/NFV-based mobile transport and computing platform supporting different verticals (e.g., automotive, e-health, e-industry) that ask for a network slice.

A more deeply description of the 5GT platform and the functional elements belonging to it are described in Section 8.1. Furthermore in Section 8.2 is presented the Scuola Superiore Sant'Anna contribution in the 5GT project.

## 8.1 Bringing the "Network Slicing" into 5G networks

The H2020 5G-TRANSFORMER project (5GT) [22] aims to transform today's mobile transport network into an SDN/NFV-based Mobile Transport and Computing Platform (MTP), which brings the "Network Slicing" paradigm into mobile transport networks by provisioning and managing MTP slices tailored to the specific needs of vertical industries. The technical approach is twofold:

- Enable vertical industries to meet their service requirements within customised MTP slices;

- Aggregate and federate transport networking and computing fabric, from the edge all the way to the core and cloud, to create and manage MTP slices throughout a federated virtualized infrastructure.

The goal of 5G-Transformer is to design, implement and demonstrate a 5G platform that addresses the aforementioned challenges. The project will demonstrate several vertical industry use cases:

- Automotive: Autonomous Cruise Control (ACC) enforcement application, Collaborative Advanced Driver Assistance Systems (ADAS) application and Remote Vehicle Interaction (RVI) application.

- eHealth: Improvement of the municipal emergency communication network and development of a new technological solution for health workers and volunteers.

- Media & Entertainment: Media applications for stadia and the Olympic Games.

5GT defines three novel building blocks that will be developed and demonstrated integrating the aforementioned three vertical industries [129]:

- The **Vertical Slicer (VS)**, that is the entry point for the verticals to request a slice. It coordinates the incoming vertical slice requests for the use of networking and computing resources. Slices are requested to the VS through a defined interface using templates (called blueprints) with simple interconnection models. It is also in charge of mapping the high-level requirements and placement constraints of the slice template into a set of one or more VNF graphs and service function chains.

- The **Service Orchestrator (SO)**, responsible for the end-to-end service orchestration. It manages the allocation and monitoring of all virtual resources to all slices. Depending on the slice requirements and network context, the SO may interact with other SOs belonging to other administrative authority domains to take decisions on the end-to-end service composition and decomposition. This can be a single or multiple administrative domains depending on resources availability and characteristics.

- The **Mobile Transport and computing Platform (MTP)** for the integration of midhaul and backhaul networks. It manages the underlying physical mobile transport network and computing infrastructure.

The work flow that involves these three components to deploy a service is represented in Fig. 8.1.



**Figure 8.1:** VS, SO and MTP vertical services instatiation high level workflow

The request is sent by a tenant to the VS. The VS translates the request into a service graph described by a NSD that follows the specification and contains also the composition of a set of VNF and the resource requirements. The VS sends the service instatiation request to the SO, that maps the service graph to an MTP network slice by means of orchestration of virtual resources to this slice. The orchestration decision for an MTP slice consists of the placement of VNFs over a virtual network as well as deciding the resources to be allocated based on the abstract view provided by the MTP. The SO sends the request to the MTP to instantiate the network slice instance. The MTP is responsible for the life-cycle of virtual resources (including networking, computing, and storage) over the underlying physical infrastructure.

### 8.1.1 Resource and Service Orchestration

The 5GT-SO is in charge of end-to-end service orchestration and federation of networking, computing and storage resources across one or multiple 5GT-MTP domains, in addition to managing the allocation of different vertical slices.

5GT-SO receives the service requirements from the VS in the shape of a NSD. To do so, the 5GT-SO performs the following steps:

1. Decides the optimal resource allocation for the whole NFVNS;

2. Decides the optimal placement of VNF;

3. Decides the optimal deployment of virtual links connecting VNF;

4. Requests federated services and/or resources when needed, in addition to some other tasks related to monitoring and management (e.g. VNF management, consistency check on requested NSDs, etc.).

Service orchestration focuses on management, instantiation, and migration of VNFs at local, edge and cloud NFVIs. The problem of mapping VNFs to (virtual) computing entities (nodes, NFVI-PoPs) and the mapping of virtual links between VNFs into (virtual) paths, depending on the granularity of abstraction offered by the 5GT-MTPs, can be tackled by different optimization strategies, namely heuristics or mixed-integer linear programming. Moreover, automatic network service management and self-configuration algorithms (e.g., failure recovery) are also required to adapt to network changes and special events triggered by the monitoring platform. In case the 5GT-SO detects that one 5GT-MTP domain alone has not

**Figure 8.2:** Overview of 5GT-SO subsystems and their interactions

enough infrastructure resources to orchestrate the required service, it interacts with other SOs via the So-So interface to compose a service across multiple federated administrative domains. In this case, 5GT-SO will dynamically discover the available administrative domains by exchanging the view with the neighboring 5GT-SOs, and negotiate with them the needed services and resources.

Figure 8.2 presents a high level overview of 5GT-SO subsystems and their interactions designed to achieve the essential 5GT-SO operation described before.

The described 5GT-SO design follows ETSI guidelines [130] and can be considered an extension of NFV-MANO. The main building blocks comprising 5GT-SO are the following:

- **NBI Exposure Layer** offers a Northbound API towards the 5GT-VS to support requests for service on-boarding, service creation, service instantiation, service modification, and service termination.

- **NFV-NS/VNF Catalogue DB/Manager** is the repository of all usable NSDs and VNFDs that can be accessed through the Catalogue Manager. A NSD is expressed in terms of chaining of VNF components and providing description of their connectivity (i.e., virtual links) and resource requirements. The NS-D/VNFD is used by the 5GT-SO in the process of NFV-NS/VNF instantiation and its life-cycle management to obtain relevant information, e.g., deploy-

ment flavors or out-scaling rules. The Catalogue Manager also takes care of the advertising of NFV-NSs for federation purpose.

- **NFVO** has the responsibility of orchestrating virtual resources across multiple domains, fulfilling the Resource Orchestration (NFVO-RO) functions, as well as of coordinating the deployment of NFV-NSs along with their life-cycle management, thus fulfilling the Network Service Orchestration (NFVO-NSO) functions. More specifically:

  - NFVO-NSO coordinates all the NFV-NS deployment operations including Authentication, Authorization and Accounting (AAA) as well as formal checks of service requests based on attributes retrieved from NSDs and VNFDs. In particular, the Composite NSO, using the algorithms implemented in the NFV-NS Orchestration Engine (NSOE), decomposes the NSDs into several segments and decides where to deploy them, i.e., whether using a local 5GT-MTP or leveraging neighbor SOs. Accordingly, the Composite NSO requests i) the Constituent NSO and then the local NFVO-RO to implement the NFVNS segment into its administrative domain; and/or ii) the federated NFVO-NSO to implement the NFV-NS segment(s) into the other administrative domains. Finally, the NFVO-NSO is responsible for the network service life-cycle management including operations such as service on-boarding, instantiation, scaling, termination, and management of the VNF forwarding graphs associated to the network services.

  - NFVO-RO maps the NFV-NS segment into a set of virtual resources through the RO Orchestration Engine (RO-OE) by deciding the placement of each VNF within the virtual infrastructure, based on specified computational, storage and networking (e.g., bandwidth) requirements. The decision is based on available virtual resources that are exposed by the 5GT-MTP via the So-Mtp Southbound Interface (SBI) or by from other domains through the So-So/East-Westbound Interface (EBI/WBI). In the latter case, the sharing of abstract views is needed to build-up a comprehensive view of resources available from different domains and is carried out by the SO-SO Resource Federation element. Then, the RO Execution Entity (RO-EE) takes care of resource provisioning by managing the coordination of correlated actions to execute/forward the allocation requests to either 5GT-MTP or to the 5GT-SO NFVO-RO of other

domains.

- **VNF Manager (VNFM)** is in charge of the life-cycle management of the VNFs deployed by the 5GT-SO using either local or remote resources (or a combination of thereof). It receives relevant VNF life-cycle events from the local NFVO and provides reconfiguration according to specified counteractions decided by the NFVO based on VNFDs (e.g., auto-scaling).

- **SO-SO Resource Advertisement** is in charge of exchanging abstract resource views (e.g., abstract topologies, computing and storage capabilities) with other domains while feeding the 5GT-SO Resource Federation entity that consolidates inputs and stores federated resources into the NFVI Resource Repository.

- **NFVI Resource Repository** stores consolidated abstract resource views received from the underlying 5GT-MTPs, either from the So-Mtp SBI or from the SO-SO Resource Federation block in case of abstract resource views received from other SOs/domains through the So-So/East- Westbound Interface (EBI/WBI).

- **NS/VNF Instance Repository** stores the instances of VNFs and NFV-NSs that have previously been instantiated.

- **SO Monitoring Service** provides the measurement reports for the 5GT-SO to support 5GT-SO monitoring management including performance monitoring and fault management, based on the collected monitoring data provided by the 5GT-MTP.

- **Service Monitoring Data Consumer** supports the life-cycle management of instantiated VNFs/NFV-NSs by collecting measurement reports from the 5GT-SO Monitoring Service and reports data to the NFVO (e.g., to trigger auto-scaling actions based on scaling rules in the NSD) and/or to the SLA Manager (e.g., to enable SLA on-line verification). Performance reports can be also used to trigger healing actions to recover from failures or service degradation. The aim is to adapt deployed services or provisioned resources while preventing service degradation due to the concurrent usage of resources from different services.

- **SLA Manager** elaborates performance reports from the Service Monitoring Data Consumer during the service life-cycle and assures that the agreed

SLAs are continuously satisfied through on-line SLA verification. In the event a requested SLA is not met, the SLA Manager may trigger scaling actions to prevent or recover from SLA violations.

Those functionalities are supposed to interwork together towards the support of a number of operations.

## 8.2 5G Network Slice Deployment

The Scuola Superiore Sant'Anna contribution in the 5GT project related to the SO is linked to the development of part of the RO-OE. The RO-OE handles the requests that are related to 5GT-MTP resources management, i.e., NFV-NS instantiation and termination, and interacts with the PA module, the core MANO and the Resource Orchestration Execution Engine WAN Infrastructure Manager (RO-EE-WIM) to accomplish the request. In particular we managed the method that translate the NSD and the VNFDs to be sent to the PA module. The PA software components expose a REST-API server for the ROE to request the calculation of resources allocations during the instantiation of NFV-NSs.

This section describes an implementation of the 5GT-VS, 5GT-SO and 5GT-MTP with the objective to experimentally demonstrate for the first time the deployment of a 5G mobile network slice through the 5GT architecture.

The 5GT-SO implementations provides the NFVO and VNFM functionality as described in [131], extended with functionality for managing PNF-containing Network Services. To do that, an additional component called Physical Network Function Manager (PNFM) has been implemented, which takes care of configuring the PNFs with the necessary parameters through RESTful HTTP messages.

The 5GT-VS prototype is a Java application based on the Spring framework, exposing its functionality through a RESTful HTTP API and a web GUI. The application is disaggregated, with its several components, i.e. slice life-cycle manager, translator, arbitrator and SouthBound (SB) drivers, communicating through Simple Message Queue Protocol (SMQP).

The role of the 5GT-MTP is played in an initial scenario directly by OpenStack, because the additional functionalities that would be provided by a full implementation of the 5GT-MTP are not needed (as further explained in Section 8.2.1). Then as described in Section 8.2.2 the 5GT-MTP has been enhanced with the Single logical Point of Contact concept introduced in ETSI IFA028 [132].

Two scenarios are analysed and deeply described in the following sections. In a first scenario, the 5GT architecture is exploited for a preliminary slice deployment utilizing edge data-center resources while in the second one the utilized resources can be also related to a core data-center and a more elaborated use of 5GT architecture is considered. In both the proposed scenarios, the openair-cn and openair-interface5g are used as mobile network software to deploy the core and the RAN respectively. The EPC and the network elements belonging to it can be deployed as individual VNF elements in a virtualised environment or can also be deployed as bundle vEPC VNF. In both scenarios, the bundle vEPC VNF deployed in an Open-Stack environment (Ocata), is considered. Regarding the RAN, both DU and CU are deployed as PNF and they utilise Option 7-1 (i.e., intra-PHY) functional split. The DU and CU are hosted in PC 6 and PC 7 of the ARNO-5G testbed respectively (see Tab. 4.1 for the PC specifications).

In both scenarios, OpenStack is deployed as a single node that includes both the controller (Ctrl) and the compute node (CN). In Openstack two networks are defined: the Openstack private network with address 10.0.0.0/24 and the Open-stack public network with address 10.10.20.0/24. The vEPC VNF ens3 interface is assigned an IP address (10.0.0.4) of the Openstack private network. A floating IP (10.10.20.112) is, then, generated from the pool of the Openstack public network addresses and it is mapped to the vEPC VNF ens3 interface address. The floating IP address allows vEPC VNF reachability. If the vEPC VNF and CU PNF are in different IP sub networks, a Virtual eXtensible LAN (VXLAN) shall be configured for the data plane interconnection.

The demo aims to show a complete network slice deployment from the service providers to the final end user. The goal of such demonstration is not related to the reliability, but to demonstrate that using the 5GT Platform a whole slice of network resources can be allocated autonomously to provide services to the end user.

## 8.2.1 5G Network Slice: Preliminary Deployment

In a preliminary 5G network slice deployment, the testbed shown in Fig. 8.3 is considered. The vEPC VNF is communicating the CU PNF, the CU PNF communicates with the DU PNF, and the User Equipment (UE) is connected to the DU PNF.

Because of Openstack configurations, the floating IP is not listed in the vEPC

**Figure 8.3:** The Preliminary 5G Network Slice Deployment Testbed

VNF IP addresses. Thus, it cannot be used in the OAI core configuration files of the vEPC VNF. Therefore, even if the vEPC VNF floating IP and the CU PNF IP (10.10.20.2) are in the same IP sub networks, the VXLAN tunnel is established between such network entities. In this way, the VXLAN interface (vxlan0) IP address (192.168.100.1) in the vEPC VNF is used in the related OAI core configuration files and for connecting it to the CU PNF, where a VXLAN interface (vxlan0) IP address (192.168.100.2) is set. At the vEPC VNF side, the configuration of VXLAN with the fixed remote IP of CU PNF is automated by startup scripts. At the CU PNF side, during the instantiation phase of NFVO life cycle event, the NFVO provides the floating IP of vEPC VNF to create the VXLAN.

The instatiation workflow is shown in Fig. 8.4. The experiment is started by requesting a mobile service at the 5GT-VS. The component translates the service request into a mobile-capable slice, and instantiates a Network Service (see Fig. 8.5) implementing such a slice through the 5GT-SO. The 5GT-SO then starts the Instantiation process: at first it requests the instantiation of a vEPC VM to Open-Stack (acting as 5GT-MTP). While booting, the vEPC VM creates one end of the VXLAN tunnel and starts the vEPC component processes (MME, HSS, S/PGW). After the instantiation of the VM is notified back to the NFVO, it starts the configuration phase. First it configures the vEPC (in this particular demo, no configuration needs to be applied) then it requests to the PNFM to configure the CU (which is represented in the Network Service as a PNF). The PNFM sends a message to the CU containing the IP of the vEPC, so that the CU can instantiate the other half of the VXLAN tunnel and establish the communication with the vEPC.

## 8.2.2 An Improved 5G Network Slice Deployment

The deployment presented in this section has been improved with respect to the environment presented in the previous section. Indeed, the novel 5GT-MTP works

**Figure 8.4:** The Preliminary 5G Network Slice Deployment Instantiation Workflow



**Figure 8.5:** Network Service Representation

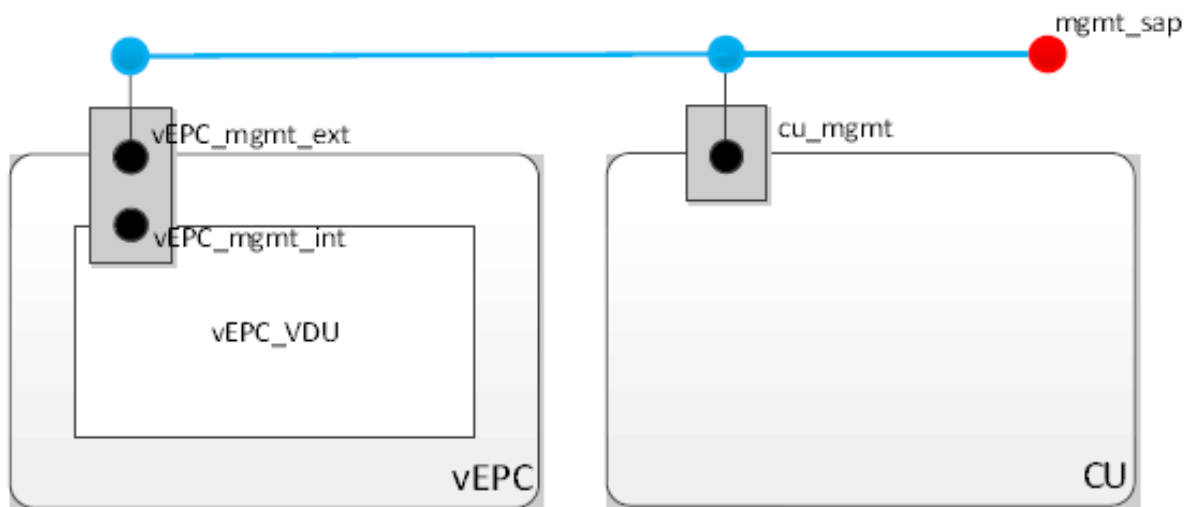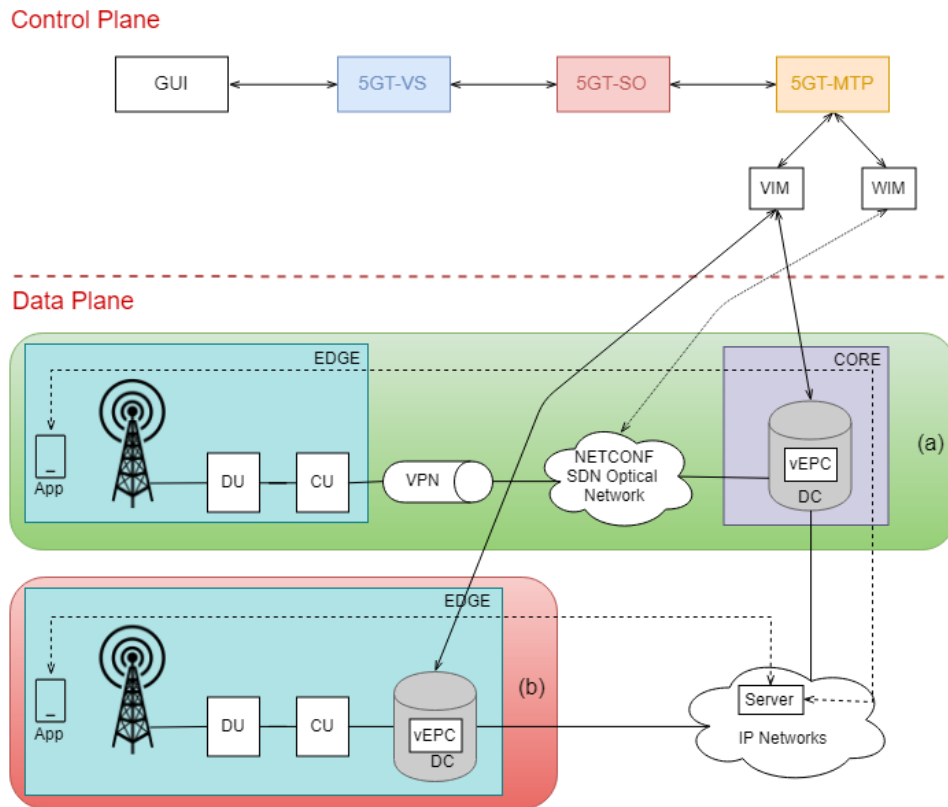as provider of NFV Infrastructure (NFVIaaS) using the Single logical Point of Contact concept introduced in ETSI IFA028. 5GT-MTP exposes one single interface to the 5GT-SO hiding the multi-VIM and multi-WIM complexity that is orchestrated and managed internally. According to ETSI IFA028, 5GT-MTP interfaces towards both the 5GT-SO and VIM/WIM are based on ETSI IFA005. 5GT-MTP exposes APIs to receive requests from the 5GT-SO and it utilises a REST client to interact by means of RESTful HTTP messages with both the Virtualized Infrastructure Manger (VIM) at the edge and the core data center and the Wide-area Infrastructure Manager (WIM) controlling the SDN optical network in the ARNO-5G testbed.

The VIM is based on OpenStack (Ocata). The WIM is based on an Open Network Operating System (ONOS) controller [133], properly extended on the Northbound and on the Southbound interfaces to operate within the 5GT architecture. Specifically, on the Northbound interface the WIM is able to receive and handle connectivity requests arriving from the 5GTMTP, using a REST API. On the Southbound interface, the ONOS-based WIM has been extended to control physical optical devices through NETCONF protocol and YANG models, as defined within the work of the ODTN (Open and Disaggregated Transport Network) project [134]. For instance, the controller is able to establish a point-to-point optical connection between two optical transponders traversing a number of ROADMs.

The Improved 5G Network Slice Deployment is able to manage two different environment: in the first one, in Fig. 8.6-a, the UEs, one Ettus B210 USPR, the DU PNF, the CU PNF are deployed at the edge premises while the vEPC, Openstack, 5GT-VS, 5GT-SO are deployed in the ARNO-5G testbed in Pisa. The two parts of the testbed (local and remote) are connected through a VPN traversing the public Internet and an SDN disaggregated optical network, including commercial devices, controlled by the WIM, based on one instance of the extended ONOS controller, in the ARNO-5G testbed.

The second scenario, in Fig. 8.6-b, includes the deployment of the vEPC within OpenStack at edge premises (in addition to what already deployed in the first scenario) and it maintains the orchestration components in the ARNO-5G testbed.

The deployment starts by requesting a mobile service through the 5GT-VS GUI. The component translates the service request into a mobile-capable slice, and instantiates a Network Service implementing such a slice through the 5GT-SO. The 5GT-SO then starts the instantiation process by interacting with the 5G-MTP: at first it requests the instantiation of a vEPC VM in OpenStack. While booting, the

**Figure 8.6:** The Improved 5G Network Slice Deployment Testbed

vEPC VM the 5GT-MTP creates one end of the VXLAN tunnel and it starts the vEPC component processes (MME, HSS, S/PGW). The successful instantiation of the VM is notified back to the NFVO and the configuration phase starts. First the NFVO VNFM configures the vEPC (in this particular demo, no configuration needs to be applied) then it requests to the PNFM to configure the CU (which is represented in the Network Service as a PNF). The PNFM sends a message to the CU containing the IP of the vEPC, so that the CU can instantiate the other half of the VXLAN tunnel and establish the communication with the vEPC. In the scenario where the slice includes part of the SDN disaggregated optical network the 5GT-SO interacts with the 5G-MTP also for establishing a lightpath between the VPN gateway and the vEPC. Once the slice is up a mobile app is capable to access the IP networks either through the remote edge vEPC (Fig. 8.6-a) or through the local edge vEPC (Fig. 8.6-b), experiencing different latencies.

## 8.3 Conclusion

In 5G, the term slicing refers, in general to the possibility for different tenants, to share the same physical networks. In this context, The H2020 5G-

TRANSFORMER project (5GT) aims at transforming today's rigid mobile transport networks into a flexible SDN/NFV-based mobile transport and computing platform supporting different verticals (e.g., automotive, e-health, e-industry) that ask for a network slice, as stated in Section 8.1.

Section 8.2.1 and Section 8.2.2 show how a novel software framework based on the 5GT architecture can deploy a mobile network slice in few minutes. Thus, the section targets the innovations demanded by the future 5G mobile networks: slice deployment specialized to vertical requirements, latency-awareness, virtualized mobile network function setup, flexibility and programmability enabled by SDN and NFV.

# 9  Closing discussion

## 9.1  Summary of the thesis

This thesis provides contributions to the state of the art regarding the architecture and performance of the LTE and 5G mobile networks, from the conceptual and the experimental point of view. The thesis covers both research and implementation aspects, especially because almost everything was developed in an experimental setup to develop and analyze a standard compliant 5G network.

The main objective of the thesis is the development of a standard compliant 5G network and all the components belonging to it. To achieve such goal, as described in Chapter 4, the openairinterface5G emulation platform has been analysed and implemented in the ARNO-5G testbed in the TeCIP Laboratory of the Scuola Superiore Sant'Anna in Pisa.

In such a network, firstly a deep analysis of its performance have been performed. This is because the need to understand if some typical 5G services and applications, such as autonomous vehicles collision application and ultra high-definition video streaming applications could be deployed exploiting the deployed 5g network. After a first phase in which the platform has been analyzed, it has been extended with new features and capabilities and adopted to setup a standard compliant 5G network. A deeply analysis and experimental measurements of latency and jitter experienced in the 5G network have been performed. As described in Chapter 5, an evaluation of how different virtualization technologies decrease the midhaul latency budget when different 5G network components are virtualised considering different gNB functional splits. In particular, Section 5.2 described the possible functional splits of the LTE-A Pro gNB between RRH and BBU. Then it analysed, through emulation, the capacity requirements for the C/VRAN midhaul. Based on the obtained results the capacity required for transporting both user and control data are limited to slightly less than one Megabit per second. Thus the limiting requirement for the midhaul is the latency which can vary from units of milliseconds, if a physical layer functional split is considered, to tens of seconds

if an RRC layer functional split is implemented. In Section 5.3 is presented an experimental evaluation of the capacity and the latency requirements of the midhaul network when the Option 7 (intra-PHY) functional split of the New Radio is implemented. As expected, the capacity requirement is independent of the traffic generated by the UE because the midhaul is carrying cell-level information. Moreover, the maximum one-way latency that can be tolerated along the midhaul is about 250 $\mu$s as specified by 3GPP. Moreover, in Section 5.4 is presented the experimental evaluation of the impact of virtualizing gNB functions on the midhaul latency budget. It also showed the maximum sustainable midhaul jitter. Results showed that by increasing the instances of CU running in the same machine the allowable midhaul latency budget decreases of some tens of microseconds due to the higher number of computations required in the same machine. Similarly, but in the order of more than fifty microseconds, it happens if the signal bandwidth increases. Moreover, if gNB functions are run in virtual machines the allowable latency budget further decreases, in the order of hundreds of microseconds, due to the higher number of computations required by the virtualization engine. Finally, the midhaul jitter evaluation showed that jitter negligibly impact the allowable latency budget. However, the allowable jitter budget is in the order of tens of microseconds in all the considered scenarios. In Section 5.5 is instead presented the experimental evaluation of the impact of virtualizing gNB functions on the midhaul latency and jitter budget when different virtualisation methods are utilized. Results showed that lighter virtualisation methods (e.g., Docker) are impacting the midhaul latency budget for Option 7-1 (i.e., intra-PHY) split less than heavier virtualisation methods (e.g., VirtualBox). However, in all the cases, the midhaul latency budget reduction depends on the considered signal bandwidth. The higher the bandwidth the higher the computations required the higher the midhaul latency budget reduction. Furthermore, the performed experimental evaluation showed that a jitter of at most 40 $\mu$s can be tolerated. Finally, in Section 5.6 an experimental analysis of the effect of virtualizing NGRAN components (e.g., CU and DU) on the maximum latency and jitter that the midhaul can support has been performed. The first set of results showed that by using heavier virtualization technologies and a higher number of physical resource blocks (i.e., channel bandwidth), the midhaul maximum latency decreases due to a heavier elaboration requested to the hardware. An empirical equation expressing the midhaul maximum latency as a linear function of the number of physical resource blocks (i.e., channel bandwidth), functional splits, and virtualization technologies confirm

the aforementioned trends. Moreover, even the midhaul jitter can be critical if it reaches values above 40 $\mu$s. A second set of results showed that if virtual DUs and CUs featuring split Option 8 are deployed, the utilization of the anti-affinity constraint is advisable to avoid large impairment in terms of maximum supported latency. A third set of results showed that by increasing the number of NG-RAN components in the same computational resource the maximum midhaul latency heavily decreases.

The next objective was the introduction in the 5G network of the so-called Multiaccess Edge Computing (MEC). Indeed as explained in Chapter 6, ETSI has published the baseline MEC standards to allow a standards-based environment for cloud applications in the network edge. Application enablement API and the MEC service APIs are the essential components in the MEC specification for a unified, standards-based environment for context-aware cloud applications. New MEC service APIs are also being developed for specific industry applications such as V2X to allow MEC better serve and add value to these applications. The 3GPP 5G system specification of Rel-15 includes native enablers for edge computing. Section 6.2 has illustrated the potential of these enablers for an integrated MEC deployment in 5G networks. The key components of this integration are the ability of MEC, as a 5G AF, to interact with the 5G system to influence the routing of the edge applications' traffic and the ability to receive notifications of relevant events, such as mobility events, in the 5G system for improved MEC deployment efficiency and end user experience. Moreover, the versatility of the 3GPP service exposure and API frameworks in principle also allows MEC to provide services to the 5G system. Section 6.3 provides an overview of automotive use cases, as introduced by 5GAA, and shows how MEC can be considered as a key technology supporting multiple services for connected AD vehicles. Moreover, this section draws attention to the value of MEC as a standardized solution for Edge Computing, especially important from automotive stakeholders' point of view while also serving other vertical market segments. In fact, a great value is associated with the standardization of edge computing technology, as open standards are the way to open the market and to ensure interoperability. In particular, from a standardization perspective, some use cases targeting fully connected cars will require the fulfillment of challenging requirements, possible only with the introduction of 5G networks. In this perspective, Section 6.4 describes a MEC-Based VRU Warning System while Section 6.5 proposes a MEC-Based orchestrator to perform a needed meticulous orchestration of different MEC services such as Advanced Driving Assistance and in-vehicle in-

fotainment services. The orchestrator proposed in Section 6.5 aims to optimize the QoE of a videostreaming service for infotainment, while guaranteeing the requested capacity to coexisting ADA services. In particular, Section 6.4 presents a complete standards-compliant design and implementation of a Vulnerable Road User Warning system. This system aims to provide an exchange of information between the road users (e.g. pedestrian, cyclist, vehicles) about the presence of nearby entities in case of dangerous situations: warnings are provided to VRUs to avoid collision with the moving vehicle and vice versa. The proposed architecture is composed by a user-side Android application that periodically sends messages containing VRU position, speed and orientation, and by a MEC-based application, named CAM server, able to process the geographical location of the VRU, predict a possible collision and, if needed, warn the involved road entities sending alert messages. A viewer showing a map with VRU and moving vehicles has been developed to help to visualize the involved entities and as demonstrator of a viewer that could be installed both on vehicles and VRUs. Moreover, in this section a preliminary performance evaluation about end-to-end latency between VRUs application and the CAM server is presented. Taking advantage of the recent advances in the area of Multi-access Edge Computing, the design and the implementation of a Connected Vehicle Service Orchestrator (CVSO) that manages heterogeneous automotive applications at the mobile edge is presented in Section 6.5. The presented system delivers QoE-optimized infotainment video services coexisting with Advanced Driving Assistance (ADA) applications, and does so in a fully standards-based way: video is delivered using DASH technologies, which are widely supported, and the proposed server-side video quality optimizations take place in a transparent way to the clients. At the same time, the proposed CSVO is deployed on top of a standards compliant MEC platform featuring a Radio Network Information Service, complying with recent ETSI standards that we have implemented. Contrary to typical receiver-driven video adaptation mechanisms, it is our CVSO that decides on the optimal video quality per user, based on the latter's channel quality characteristics which are available via the RNIS. This involves solving an ILP. Via testbed experiments over a MEC-capable LTE network, it has been shown the proposed server-assisted video adaptation scheme to improve on user experience, while also showing that the ILP solution can be derived in reasonable time even for large numbers of users, thus verifying the proposed system's feasibility and suitability for MEC deployment. This objective is mainly realized thanks to the collaboration with EURECOM in Sophia-Antipolis during the period abroad.

The 5G network architecture, will be heavily based on virtual network functions( VNFs). Thus, in a scenario where network functions are virtualized, both hardware and software failures assume the same importance, and their reliability shall be guaranteed. Similarly, reliability at service chain level is important to assure proper service availability features to application service platforms deployed by verticals. Therefore a study of the reliability in 5G networks has been also performed in the conclusion of my Ph.D. Different protection mechanisms have been proposed and experimentally analysed for both 5G core and access networks. Chapter 7 focuses on this topics and therefore described a testbed conceived to test reliability mechanisms applied to both backhaul and midhaul in NG RAN. In Section 7.1, a resilience scheme for the NG Core is proposed. It is shown the performance of resilience schemes for virtualised mobile network functions. In particular, the Section 7.1 focused on evaluating the time required to regain UE connectivity when a hot backup vEPC is deployed close to or away from a working vEPC. The study exploited the remote access to federated testbeds. Section 7.2 described a testbed conceived to test reliability mechanisms applied to midhaul in NG RAN. The use case in this study was that the NG RAN functional split option 7-1 was provided by OAI running on a number of Linux servers. The midhaul were implemented using either a simple Ethernet cable or the two layer Ethernet-over-DWDM PROnet. Two reliability mechanisms were discussed and applied to the midhaul to improve its reliability, defined as the ability of the network not to disconnect the UE from the NG RAN in the presence of fiber span outage. The first mechanism makes use of 1:1 protection at the Ethernet layer combined with an optical circuit restoration at the DWDM layer. The second mechanism makes use of 1 + 1 protection at the Ethernet layer combined with an optical circuit restoration at the DWDM layer. These two combined protection and restoration mechanisms (1:1 + R and 1 + 1 + R) are coordinated by the software defined network (SDN) PROnet orchestrator to overcome any single fiber span outage, which may be followed by other subsequent outages. More specifically, the implementation of the 1 + 1 protection mechanism required the development of a new kernel software module, which runs in Open vSwitch (OVS) and achieves zero packet loss even during the occurrence of a fiber outage. Experimentally, it has been determined that in the midhaul, the 1:1 protection mechanism was failed to keep the UE connected to the C-RAN, due to the prolonged traffic disruption and the resulting burst of lost data packets. However, by applying the 1 + 1 protection mechanism to the midhaul, the UE remained connected to the C-RAN even after the fiber outage

occurrence. Finally, Section 7.3 proposed a two-step scheme for recovering virtualized distributed unit (vDU) and virtualized central unit (vCU) connectivity upon midhaul failure. In particular, a single soft failure of a lightpath interconnecting vCU and vDU is considered. The main novelty of the proposed scheme consists of complementing lightpath adaptation with the possibility of flexibly changing the vCU and vDU functional split. This approach allows the increase of the number of the recovered connections when spare network resources are scarce. Simulation results showed that the proposed scheme allows us to recover almost all the midhaul vDU-vCU connections along the same path where working lightpaths were routed, and it largely overcomes the performance of lightpath rerouting, especially when the network load is high. Experimental results showed that the recovery time experienced when functional split change is triggered is in the order of tens of seconds.

Last but not least, this work also pursued a really challenging goal: we were able to demonstrate a 5G network slice deployment in the ARNO testbed by using the 5G-TRANSFORMER architecture and offer a mobile/edge connectivity service with virtualized functions. In 5G, the term slicing refers, in general to the possibility for different tenants, to share the same physical networks. In this context, The H2020 5G TRANSFORMER project (5GT) aims at transforming today's rigid mobile transport networks into a flexible SDN/NFV-based mobile transport and computing platform supporting different verticals (e.g., automotive, e-health, e-industry) that ask for a network slice, as stated in Section 8.1. Section 8.2.1 and Section 8.2.2 show how a novel software framework based on the 5GT architecture can deploy a mobile network slice in few minutes. Thus, the section targets the innovations demanded by future mobile networks: slice deployment specialized to vertical requirements, latency-awareness, virtualized mobile network function setup, flexibility and programmability enabled by SDN and NFV.

## 9.2 Ph.D. Outcomes

**Journal Papers**

[1] Kondepu, K.; Sgambelluri, A.; Sambo, N.; **Giannone, F.**; Castoldi, P.; Valcarenghi, L.; **"Orchestrating lightpath recovery and flexible functional split to preserve virtualized RAN connectivity"**, IEEE/OSA Journal of Optical Communications and Networking, 10, 11,

843-851, 2018/11/21, IEEE.

[2] **Giannone, F.**; Kondepu, K.; Gupta, H.; Civerchia, F.; Castoldi, P.; Franklin, A. Antony; Valcarenghi, L.; **"Impact of Virtualization Technologies on Virtualized RAN Midhaul Latency Budget: A Quantitative Experimental Evaluation"**, IEEE Communications Letters, 23, 4, 604-607, 2019/2/14, IEEE.

[3] Ramanathan, S.; Tacca, M.; Razo, M.; Mirkhanzadeh, B.; Kondepu, K.; **Giannone, F.**; Valcarenghi, L.; Fumagalli, A.; **"A programmable optical network testbed in support of C-RAN: a reliability study"**, Photonic Network Communications, 37, 3, 311-321,2019/6/15, Springer, US.

**Conference Papers**

[1] Valcarenghi, L; Kondepu, K; **Giannone, F.**; Castoldi, P; "Requirements for 5G fronthaul", 2016 18th IEEE International Conference on Transparent Optical Networks (ICTON), 1-5, 2016/7/10, IEEE;

[2] Marotta A., Kondepu K., **Giannone, F.**; Doddikrinda S., Cassioli D., Antonelli C., Valcarenghi L., Castoldi P., **"Performance evaluation of CoMP coordinated scheduling over different backhaul infrastructures: A real use case scenario"**, 2016 IEEE International Conference on the Science of Electrical Engineering (ICSEE), 1-5, 2016/11/16, IEEE;

[3] Marotta, A.; **Giannone, F.**; Kondepu, K.; Cassioli, D.; Antonelli, C.; Valcarenghi, L.; Castoldi, P., **"Reducing CoMP control message delay in PON backhauled 5G networks"**, European Wireless 2017; 23th European Wireless Conference, 1-5, 2017/5/17, VDE;

[4] Marotta, A.; Kondepu, K.; **Giannone, F.**; Cassioli, D.; Antonelli, C.; Valcarenghi, L.; Castoldi, P.; **"Impact of CoMP VNF placement on 5G coordinated scheduling performance"**, 2017 European Conference on Networks and Communications (EuCNC), 1-6, 2017/6/12, IEEE;

[5] Valcarenghi, L.; **Giannone, F.**; Manicone, D; Castoldi, P.; **"Virtualized eNB latency limits"**, 19th International Conference on Transparent Optical Networks (ICTON), 1-4, 2017/7/2, IEEE;

[6] Gupta, H.; Manicone, D.; **Giannone, F.**; Kondepu, K.; Franklin, A.; Castoldi, P.; Valcarenghi, L.; **"How much is fronthaul latency budget impacted by RAN virtualisation ?"**, 2017 IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN) - Workshop on Federated Testbeds for NFV/SDN/5G: Experiences and Feedbacks,315-320,2017/11/7, IEEE.

[7] Kondepu, K.; Giorgetti, A.; **Giannone, F.**; Marotta, A.; Cugini, F.; Castoldi, P.; Valcarenghi, L.; "SDN-enabled Latency-Guaranteed Dual Connectivity in 5G RAN", 2017 Asia Communications and Photonics Conference (ACP), 1-3, 2017/11/10, IEEE.

[8] **Giannone, F.**; Gupta, H.; Kondepu, K.; Manicone, D.; Franklin, A.; Castoldi, P.; Valcarenghi, L.; **"Impact of RAN virtualization on fronthaul latency budget: An experimental evaluation**,2017 IEEE Globecom Workshops (GC Wkshps), 1-5, 2017/12/4, IEEE.

[9] Kondepu, K.; Sambo, N.; **Giannone, F.**; Castoldi, P.; Valcarenghi, L.; **"Orchestrating lightpath adaptation and flexible functional split to recover virtualized RAN connectivity"**, 2018 Optical Fiber Communications Conference and Exposition (OFC), 1-3, 2018/3/11, IEEE.

[10] Ponzini, F.; Kondepu, K.; **Giannone, F.**; Castoldi, P.; Valcarenghi, L.; **"Optical Access Network Solutions for 5G Fronthaul"**, 2018 20th International Conference on Transparent Optical Networks (ICTON), 1-5, 2018/7/1,IEEE.

[11] Civerchia, F.; Kondepu, K.; **Giannone, F.**; Doddikrinda, S.; Castoldi, P.; Valcarenghi, L.; **"Encapsulation techniques and traffic characterisation of an Ethernet-based 5G fronthaul"**, 2018 20th International Conference on Transparent Optical Networks (ICTON), 1-5, 2018/7/1, IEEE.

[12] Napolitano, A.; **Giannone, F.**; Civerchia, F.; Kondepu, K.; Cecchetti, G.; Ruscelli, A.; Valcarenghi, L.; Castoldi, P.; **"Italian 5G Trials: A Vertical View"**, 2018 IEEE 4th International Forum on Research and Technology for Society and Industry (RTSI), 1-5, 2018/9/10, IEEE.

[13] Capitani, M.; **Giannone, F.**; Fichera, S.; Sgambelluri, A.; Kondepu, K.; Kraja, E.; Martini, B.; Landi, G.; Valcarenghi, L.; **"Experimental Demonstration of a 5G Network Slice Deployment Through the 5G-Transformer Architecture"**, 2018 European Conference on Optical Communication (ECOC), 1-3, 2018/9/23, IEEE.

[14] Sgambelluri, A.; Capitani, M.; Fichera, S.; Kondepu, K.; Giorgetti, A.; **Giannone, F.**; Martini, B.; Ubaldi, F.; Iovanna, P.; Landi, G.; **"Experimental Demonstration of a 5G Network Slice Deployment Exploiting Edge or Cloud Data-Centers"**, 2019 Optical Fiber Communications Conference and Exhibition (OFC), 1-3,2019/3/3,IEEE.

[15] Kondepu, K.; **Giannone, F.**; Vural, S.; Riemer, B.; Castoldi, P.; Valcarenghi, L.; **"Experimental Demonstration of 5G Virtual EPC Recovery in Federated Testbeds"**, 2019/4/8, IFIP/IEEE Symposium on Integrated Network and Service Management (IM), 712-713, 2019/4/8, IEEE.

[16] Napolitano, A.; Cecchetti, G.; **Giannone, F.**; Ruscelli, A.; Civerchia, F.; Kondepu, K.; Valcarenghi, L.; Castoldi, P.; **"Implementation of a MEC-based Vulnerable Road User Warning System"**, 2019 AEIT International Conference of Electrical and Electronic Technologies for Automotive (AEIT AUTOMOTIVE), 1-6,2019/7/2, IEEE.

# Bibliography

[1] T. Rappaport, *Wireless communications: Principles and practice*. Prentice Hall, 1996.

[2] H. Lacohée, N. Wakeford, and I. Pearson, "A social history of the mobile telephone with a view of its future," *BT Technology Journal*, vol. 21, no. 3, pp. 203–211, Jul 2003. [Online]. Available: https://doi.org/10.1023/A:1025187821567

[3] T. Dunnewijk and S. Hultén, "A brief history of mobile communication in europe," *Telematics and Informatics*, vol. 24, pp. 164–179, 08 2007.

[4] K. W. Richardson, "Umts overview," *Electronics Communication Engineering Journal*, vol. 12, no. 3, pp. 93–100, June 2000.

[5] H. Holma, A. Toskala, and J. Reunanen, *Further Outlook for LTE Evolution and 5G*. Wiley, 2015. [Online]. Available: https://ieeexplore.ieee.org/document/8042886

[6] O. Oshin, M. Luka, and P. Atayero, *From 3GPP LTE to 5G: An Evolution*. Springer, 2016.

[7] K. Sridhar, "Introduction to Evolved Packet Core," Alcatel - Lucent, White Paper, Aug. 2012.

[8] E. Dahlman, S. Parkvall, and J. Skold, *4G LTE/LTE-advanced for mobile broadband*. Elsevier, 2013.

[9] H. Holma and A. Toskala, *LTE for UMTS - OFDMA and SC - FDMA Based Radio Access*. Wiley, 2009.

[10] K. Sridhar, "3GPP Long Term Evolution: System Overview, Product Development, and Test Challenges," Agilent Technologies, Application Note, Aug. 2009.

[11] "3GPP Mobile Broadband Innovation Path to 4G: Release 9, Release 10 and Beyond: HSPA+, SAE/LTE and LTE-Advanced," 3G Americas, White Paper, Oct. 2009.

[12] A. A. Atayero, M. K. Luka, M. K. Orya, and J. O. Iruemi, "3GPP Long Term Evolution: Architecture, Protocols and Interfaces," *International Journal of Information and Communication Technology Research*, vol. 1, no. 7, pp. 306–310, Nov. 2011.

[13] O. I. Oshin and A. A. Atayero, "3GPP LTE: An Overview," in *World Congress on Engineering*, Toronto, Ontario, Canada, Jul. 2015.

[14] M. Elkashlan, T. Q. Duong, and H. Chen, "Millimeter-wave communications for 5G - Part 2: applications," *IEEE Communications Magazine*, vol. 53, no. 1, pp. 166–167, January 2015.

[15] A. Osseiran *et al.*, "Scenarios for 5G mobile and wireless communications: the vision of the METIS project," *IEEE Communications Magazine*, vol. 52, no. 5, pp. 26–35, May 2014.

[16] G. Brown, "Designing Cloud-Native 5G Core Networks," Heavy Reading, White Paper produced for Nokia, Feb. 2017.

[17] ——, "New Transport Network Architectures for 5G RAN," Heavy Reading, White Paper produced for Fujitsu, Feb. 2017.

[18] A. Galis. (2017) "Network Slicing Terms and Systems". [Online]. Available: https://datatracker.ietf.org/meeting/99/materials/slides-99-netslicing-alex-galis-netslicing-terms-and-systems

[19] NGMN P1 WS1 E2E Architecture team, "Description of Network Slicing Concept," NGMN, Final Deliverable (approved) 23.007, Sept. 2017, v1.0.8 (draft).

[20] 3GPP, "3rd Generation Partnership Project; Technical Specification Group Services and System Aspects; Telecommunication management; Study on management and orchestration of network slicing for next generation network," 3rd Generation Partnership Project (3GPP), Technical Report (TR) 28.801, Jan. 2018, version 15.1.0 Release 15.

[21] ETSI GGR NFV-EVE 012, "Network Functions Virtualisation (NFV) Release 3; Evolution and Ecosystem; Report on Network Slicing Support with ETSI NFV Architecture Framework," ETSI, Group Report, Dec. 2017, v3.1.1.

[22] 5G-Transformer. 5G Mobile Transport Platform for Verticals. [Online]. Available: http://5g-transformer.eu/

[23] EURECOM. Graduate School and Research Center in Digital Sciences. [Online]. Available: http://www.eurecom.fr

[24] OSA. Openairinterface - 5g software alliance for democratising wireless innovation. [Online]. Available: https://www.openairinterface.org/

[25] TeCIP. (2001) Institute of Communication, Information and Perception Technologies. [Online]. Available: https://www.santannapisa.it/en/institute/tecip/tecip-institute

[26] The ARNO Testbed. [Online]. Available: http://arnotestbed.santannapisa.it/

[27] jfed. [Online]. Available: http://doc.fed4fire.eu/getanaccount.html

[28] P. Rost, C. J. Bernardos, A. D. Domenico, M. D. Girolamo, M. Lalam, A. Maeder, D. Sabella, and D. WÃ¼bben, "Cloud technologies for flexible 5G radio access networks," *IEEE Communications Magazine*, vol. 52, no. 5, pp. 68–76, May 2014.

[29] T. Pfeiffer, "Next generation mobile fronthaul and midhaul architectures," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 7, no. 11, pp. B38–B45, November 2015.

[30] 3GPP, "LTE; Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description; Stage 2," 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 36.300, 04 2016, version 13.3.0 Release 13.

[31] U. Dötsch, M. Doll, H. Mayer, F. Schaich, J. Segel, and P. Sehier, "Quantitative analysis of split base station processing and determination of advantageous architectures for LTE," *Bell Labs Technical Journal*, vol. 18, no. 1, pp. 105–128, June 2013.

[32] A. Maeder, M. Lalam, A. De Domenico, E. Pateromichelakis, D. Wübben, J. Bartelt, R. Fritzsche, and P. Rost, "Towards a flexible functional split for cloud-RAN networks," in *2014 European Conference on Networks and Communications (EuCNC)*, June 2014, pp. 1–5.

[33] D. Wübben, P. Rost, J. S. Bartelt, M. Lalam, V. Savin, M. Gorgoglione, A. Dekorsy, and G. Fettweis, "Benefits and Impact of Cloud Computing on 5G Signal Processing: Flexible centralization through cloud-RAN," *IEEE Signal Processing Magazine*, vol. 31, no. 6, pp. 35–44, Nov 2014.

[34] N. Nikaein, M. K. Marina, S. Manickam, A. Dawson, R. Knopp, and C. Bonnet, "OpenAirInterface: A Flexible Platform for 5G Research," *SIGCOMM Comput. Commun. Rev.*, vol. 44, no. 5, pp. 33–38, Oct. 2014. [Online]. Available: http://doi.acm.org/10.1145/2677046.2677053

[35] 3GPP, "Technical Specification Group Radio Access Network; Evolved Universal Terrestrial Radio Access (E-UTRA) Medium Access Control (MAC) protocol specification," 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 36.300, 12 2007, version 8.0.0 Release 8.

[36] J. T. J. Penttinen, *The Telecommunications Handbook: Engineering Guidelines for Fixed, Mobile and Satellite Systems*. John Wiley & Sons, 2015.

[37] C. Chang, N. Nikaein, and T. Spyropoulos, "Impact of Packetization and Scheduling on C-RAN Fronthaul Performance," in *2016 IEEE Global Communications Conference (GLOBECOM)*, Dec 2016, pp. 1–7.

[38] N. Nikaein, "Processing Radio Access Network Functions in the Cloud: Critical Issues and Modeling," in *Proceedings of the 6th International Workshop on Mobile Cloud Computing and Services*, ser. MCS '15, 2015, pp. 36–43.

[39] *Mobile Edge Computing (MEC); Mobile Edge Management; Part 1: System host and platform management*, ETSI Group Specification MEC 010-1 V1.1.1, Jul. 2017.

[40] *Mobile Edge Computing (MEC); Mobile Edge Management; Part 2: Application lifecycle, rules and requirements management*, ETSI Group Specification MEC 010-2 V1.1.1, Jul. 2017.

[41] *Mobile Edge Computing (MEC);Mobile Edge Platform Application Enablement*, ETSI Group Specification MEC 011 V1.1.1, Jul. 2017.

[42] *Mobile Edge Computing (MEC);Radio Network Information API*, ETSI Group Specification MEC 012 V1.1.1, Jul. 2017.

[43] *Mobile Edge Computing (MEC);Location API*, ETSI Group Specification MEC 013 V1.1.1, Jul. 2017.

[44] *Mobile Edge Computing (MEC);UE Identity API*, ETSI Group Specification MEC 014 V1.1.1, Feb. 2018.

[45] *Mobile Edge Computing (MEC);Bandwidth Management API*, ETSI Group Specification MEC 015 V1.1.1, Oct. 2017.

[46] *Mobile Edge Computing (MEC);UE Application Interface*, ETSI Group Specification MEC 016 V1.1.1, Sept. 2017.

[47] 3GPP, "Technical Specification Group Services and System Aspects; System Architecture for the 5G System; Stage 2," 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 23.501, Mar. 2018, version 15.1.0 Release 15.

[48] F. Giust, G. Verin, K. Antevski, J. Chou, Y. Fang, W. Featherstone, F. Fontes, D. Frydman, A. M. A. Li, D. Purkayastha, D. Sabella, C. Wehner, K. W. Wen, and Z. Zhou, "MEC deployments in 4G and evolution towards 5G," ETSI, White Paper, Feb. 2018.

[49] *Mobile Edge Computing (MEC); Framework and Reference Architecture*, ETSI Group Specification MEC 003 V1.1.1, Mar. 2016.

[50] 5GAA, "Use Case and KPI requirements: Prioritization and Timeline; Interims Status," 5G Automotive Association, Technical Report T-170215, 2017, v1.0.

[51] DOT HS 811 731, "Description of Light-Vehicle Pre-Crash Scenarios for Safety Applications Based On Vehicle-to-Vehicle Communications," National Highway Traffic Safety Administration, Department of Transportation, May 2013.

[52] DOT HS 812 312, "Crash Avoidance Needs and Countermeasure Profiles for Safety Applications Based on Light-Vehicle-to-Pedestrian Communications," National Highway Traffic Safety Administration, Department of Transportation, Aug. 2016.

[53] 3GPP, "Technical Specification Group Services and System Aspects; Study on LTE support for Vehicle to Everything (V2X) services," 3rd Generation Partnership Project (3GPP), Technical Report (TR) 22.885, Dec. 2015, version 14.0.0 Release 14.

[54] A. Voronov, "ETSI ITS G5 GeoNetworking Stack," https://github.com/alexvoronov/geonetworking/tree/master/camdenm.

[55] Federal Communications Commission, "Chapter 73.208," https://www.gpo.gov/fdsys/pkg/CFR-2005-title47-vol4/pdf/CFR-2005-title47-vol4-sec73-208.pdf.

[56] I. Sodagar, "The MPEG-DASH standard for multimedia streaming over the internet," *IEEE MultiMedia*, vol. 18, no. 4, pp. 62–67, 2011.

[57] X. Foukas, N. Nikaein, M. M. Kassem, M. K. Marina, and K. P. Kontovasilis, "Flexran: A flexible and programmable platform for software-defined radio access networks," in *Proc. ACM CoNEXT*, 2016.

[58] S. Arora, P. A. Frangoudis, and A. Ksentini, "Exposing radio network information in a MEC-in-NFV environment: the RNISaaS concept," EURECOM, Research Report RR-19-339, Feb. 2019.

[59] *Information Technology - Dynamic Adaptive Streaming over HTPP (DASH) - Part 1: Media Presentation Description and Segment Formats*, ISO/IEC Standard 23 009.1:2014, May 2014.

[60] *Dynamic adaptive streaming over HTTP (DASH) – Part 5: Server and network assisted DASH (SAND)*, ISO/IEC Standard 23009-5:2017, Feb. 2017.

[61] E. Thomas, M. O. van Deventer, T. Stockhammer, A. C. Begen, M. L. Champel, and O. Oyman, "Applications and deployments of server and network assisted DASH (SAND)," in *Proc IBC*, 2016.

[62] *3rd Generation Partnership Project; Technical Specification Group Radio Access Network; Evolved Universal Terrestrial Radio Access (E-UTRA); Physical layer procedures*, 3GPP, Technical Specification (TS), Sept. 2018, Version 13.11.0.

[63] *Subjective Audiovisual Quality Assessment Method for Multimedia Applications*, ITU-T Recommendation, P.911, Dec. 1998.

[64] P. Schmitt, B. Landais, and F. Y. Yang, "Control and User Plane Separation of EPC nodes (CUPS)," 3GPP, Tech. Rep., Jul. 2018. [Online]. Available: http://www.3gpp.org/cups

[65] "Elephants Dream," https://orange.blender.org/.

[66] "FFmpeg tool," https://trac.ffmpeg.org/.

[67] "GPAC Multimedia Open Source Project," https://gpac.wp.imt.fr/.

[68] S. Lederer. Optimal Adaptive Streaming Formats MPEG-DASH & HLS Segment Length. [Online]. Available: https://bitmovin.com/mpeg-dash-hls-segment-length/

[69] J. Allnatt, *Transmitted-Picture Assessment*. Wiley, 1983.

[70] *Parametric non-intrusive assessment of audiovisual media streaming quality. Amendment 2: New Appendix III - Use of ITU-T P.1201 for non-adaptive, progressive download type media streaming*, ITU-T Recommendation, P.1201 Standard, Dec. 2013.

[71] M. N. Garcia, R. Schleicher, and A. Raake, "Impairment-Factor-Based Audiovisual Quality Model for IPTV: Influence of Video Resolution, Degradation Type, and Content Type," *EURASIP Journal on Image and Video Processing"*, vol. 2011, no. 1, Mar. 2011.

[72] *The E-model: A computational model for use in transmission planning*, ITU-T Recommendation, G.107 Standard, 2013.

[73] Y. Li, P. A. Frangoudis, Y. Hadjadj-Aoul, and P. Bertin, "A Mobile Edge Computing-assisted video delivery architecture for wireless heterogeneous networks," in *Proc. IEEE ISCC*, 2017.

[74] G. Rubino, "Quantifying the Quality of Audio and Video Transmissions over the Internet: The PSQA Approach," in *Communication Networks & Computer Systems*, J. A. Barria, Ed. Imperial College Press, 2005.

[75] K. D. Singh, Y. Hadjadj-Aoul, and G. Rubino, "Quality of Experience estimation for Adaptive HTTP/TCP video streaming using H.264/AVC," in *Proc. IEEE CCNC*, 2012.

[76] T. Groléat, M. Sbai, S. Vaton, Y. Hadjadj-Aoul, K. Singh, and S. Moteau, *Advances on monitoring primitives and integration in the VIPEER prototype*, Dec. 2012.

[77] "IBM ILOG CPLEX Optimization Studio," https://www.ibm.com/products/ilog-cplex-optimization-studio.

[78] M. Sauter, *From Gsm To Lte-Advanced Pro And 5G: An Introduction To Mobile Networks And Mobile Broadband, 3Rd Edition*.    Elsevier, Wiley.

[79] V. Nguyen, A. Brunstrom, K. Grinnemo, and J. Taheri, "SDN/NFV-Based Mobile Packet Core Network Architectures: A Survey," *IEEE Communications Surveys Tutorials*, vol. 19, no. 3, pp. 1567–1602, thirdquarter 2017.

[80] M. Gharbaoui, C. Contoli, G. Davoli, G. Cuffaro, B. Martini, F. Paganelli, W. Cerroni, P. Cappanera, and P. Castoldi, "Experimenting latency-aware and reliable service chaining in Next Generation Internet testbed facility," in *2018 IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN)*, Nov 2018, pp. 1–4.

[81] K. Antevski, J. Martín-Pérez, N. Molner, C. Chiasserini, F. Malandrino, P. A. Frangoudis, A. Ksentini, X. Li, J. SalvatLozano, R. Martínez, I. Pascual, J. Mangues-Bafalluy, J. Baranda, B. Martini, and M. Gharbaoui, "Resource Orchestration of 5G Transport Networks for Vertical Industries," *CoRR*, vol. abs/1807.10430, 2018. [Online]. Available: http://arxiv.org/abs/1807.10430

[82] 3GPP, "Technical Specification Group Core Network and Terminals; Restoration procedures," 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 23.007, Dec. 2017, version 14.4.0 Release 14.

[83] *Network Functions Virtualisation (NFV); Reliability; Report on Scalable Architectures for Reliability Management*, Group Specification, Sept. 2015.

[84] F. F. Moghaddam, A. Gherbi, and Y. Lemieux, "Self-Healing Redundancy for OpenStack Applications through Fault-Tolerant Multi-Agent Task Scheduling," in *2016 IEEE International Conference on Cloud Computing Technology and Science (CloudCom)*, Dec 2016, pp. 572–577.

[85] C. Colman-Meixner, G. B. Figueiredo, M. Fiorani, M. Tornatore, and B. Mukherjee, "Resilient cloud network mapping with virtualized BBU

placement for cloud-RAN," in *2016 IEEE International Conference on Advanced Networks and Telecommunications Systems (ANTS)*, Nov 2016, pp. 1–3.

[86] K. Kondepu, A. Sgambelluri, N. Sambo, F. Giannone, P. Castoldi, and L. Valcarenghi, "Orchestrating lightpath recovery and flexible functional split to preserve virtualized RAN connectivity," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 10, no. 11, pp. 843–851, Nov 2018.

[87] SoftFIRE. Software Defined Networks and Network Function Virtualization Testbed within FIRE+. [Online]. Available: https://www.softfire.eu/

[88] Software Documentation . SoftFIRE Middleware. [Online]. Available: http://docs.softfire.eu/softfire-middleware/

[89] OpenStack. Build the future of Open Infrastructure. [Online]. Available: https://www.openstack.org/

[90] OPEN BATON. An extensible and customizable NFV MANO-compliant framework. [Online]. Available: https://openbaton.github.io/

[91] K. Kondepu, S. Ramanathan, M. Tacca, M. Razo, B. Mirkhanzadeh, F. Giannone, L. Valcarenghi, and A. Fumagalli, "Experimental demonstration of a packet-based protection for seamlessly recovering from a multi-layer metro network fronthaul failure," in *IEEE INFOCOM 2019 - IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*, April 2019, pp. 906–908.

[92] *40-Gigabit-Capable Passive Optical Networks (NG-PON2)*, ITU-T Recommendation, G.989 Standard, Oct. 2015.

[93] A. Marotta, K. Kondepu, D. Cassioli, C. Antonelli, L. M. Correia, and L. Valcarenghi, "Software Defined 5G Converged Access as a viable Techno-Economic Solution," in *2018 Optical Fiber Communications Conference and Exposition (OFC)*, March 2018, pp. 1–3.

[94] G. Talli, S. Porto, D. Carey, N. Brandonisio, P. Ossieur, P. Townsend, R. Bonk, T. Pfeiffer, F. Slyne, S. McGettrick, C. Blümm, M. Ruffini, A. Hill, D. Payne, and N. Parsons, "Technologies and architectures to enable SDN in converged 5G/optical access networks," in *2017 International Conference on Optical Network Design and Modeling (ONDM)*, May 2017, pp. 1–6.

[95] X. Wang, C. Cavdar, L. Wang, M. Tornatore, H. S. Chung, H. H. Lee, S. M. Park, and B. Mukherjee, "Virtualized Cloud Radio Access Network for 5G Transport," *IEEE Communications Magazine*, vol. 55, no. 9, pp. 202–209, Sep. 2017.

[96] J. Rak, *Resilient Routing in Communication Networks*. Springer, 2015.

[97] S. Ramamurthy, L. Sahasrabuddhe, and B. Mukherjee, "Survivable WDM mesh networks," *Journal of Lightwave Technology*, vol. 21, no. 4, pp. 870–883, April 2003.

[98] B. Mirkhanzadeh, A. Shakeri, C. Shao, M. Razo, M. Tacca, G. Galimberti, G. Martinelli, M. Cardani, and A. Fumagalli, "An SDN-enabled multi-layer protection and restoration mechanism," *Optical Switching and Networking*, vol. 30, pp. 23–32, Nov. 2018.

[99] Open Networking Foundation. [Online]. Available: www.opennetworking. org

[100] Open vSwitch. [Online]. Available: www.openvswitch.org

[101] D. Hicks, C. Malina-Maxwell, M. Razo, M. Tacca, A. Fumagalli, and D. Nguyen, "PROnet: A programmable optical network prototype," in *2016 18th International Conference on Transparent Optical Networks (ICTON)*, July 2016, pp. 1–4.

[102] OpenFlow Switch Specification. [Online]. Available: www.opennetworking. org/wpcontent/uploads/2013/04/openflow-spec-v1.3.1.pdf

[103] D. Chitimalla, K. Kondepu, L. Valcarenghi, M. Tornatore, and B. Mukherjee, "5G fronthaul-latency and jitter studies of CPRI over Ethernet," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 9, no. 2, pp. 172–182, Feb 2017.

[104] F. Giannone, H. Gupta, K. Kondepu, D. Manicone, A. Franklin, P. Castoldi, and L. Valcarenghi, "Impact of RAN Virtualization on Fronthaul Latency Budget: An Experimental Evaluation," in *2017 IEEE Globecom Workshops (GC Wkshps)*, Dec 2017, pp. 1–5.

[105] *GSTR-TN5G Transport network*, ITU-T Technical Report Telecommunication Standardization Sector of ITU, 2018.

[106] J. Zou, C. Wagner, and M. Eiselt, "Optical Fronthauling for 5G Mobile: A Perspective of Passive Metro WDM Technology," in *Optical Fiber Communication Conference*. Optical Society of America, 2017, p. W4C.2.

[107] G. N. Liu, L. Zhang, T. Zuo, and Q. Zhang, "IM/DD Transmission Techniques for Emerging 5G Fronthaul, DCI, and Metro Applications," *Journal of Lightwave Technology*, vol. 36, no. 2, pp. 560–567, Jan 2018.

[108] V. Passas, V. Miliotis, N. Makris, T. Korakis, and L. Tassiulas, "Paris Metro Pricing for 5G HetNets," in *2016 IEEE Global Communications Conference (GLOBECOM)*, Dec 2016, pp. 1–6.

[109] Y. Pointurier, "Design of low-margin optical networks," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 9, no. 1, pp. A9–A17, Jan 2017.

[110] N. Sambo, F. Cugini, A. Sgambelluri, and P. Castoldi, "Monitoring Plane Architecture and OAM Handler," *Journal of Lightwave Technology*, vol. 34, no. 8, pp. 1939–1945, Apr. 2016.

[111] A. Dupas, P. Layec, E. Dutisseuil, S. Bigo, S. Belotti, S. Misto, S. Annoni, Y. Yan, E. Hugues-Salas, G. Zervas, and D. Simeonidou, "Hitless 100 Gbit/s OTN bandwidth variable transmitter for software-defined networks," in *2016 Optical Fiber Communications Conference and Exhibition (OFC)*, March 2016, pp. 1–3.

[112] P. Layec, A. Dupas, A. Bisson, and S. Bigo, "QoS-aware protection in flex-grid optical networks," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 10, no. 1, pp. A43–A50, Jan 2018.

[113] K. Kondepu, N. Sambo, F. Giannone, P. Castoldi, and L. Valcarenghi, "Orchestrating Lightpath Adaptation and Flexible Functional Split to Recover Virtualized RAN Connectivity," in *2018 Optical Fiber Communications Conference and Exposition (OFC)*, March 2018, pp. 1–3.

[114] M. R. Rahman and R. Boutaba, "SVNE: Survivable Virtual Network Embedding Algorithms for Network Virtualization," *IEEE Transactions on Network and Service Management*, vol. 10, no. 2, pp. 105–118, June 2013.

[115] F. Gu, H. Alazemi, A. Rayes, and N. Ghani, "Survivable Cloud Networking Services," in *2013 International Conference on Computing, Networking and Communications (ICNC)*, Jan 2013, pp. 1016–1020.

[116] L. Valcarenghi, F. Cugini, F. Paolucci, and P. Castoldi, "Quality-of-service-aware fault tolerance for grid-enabled applications," *Optical Switching and Networking*, vol. 5, pp. 150–158, Jun. 2008.

[117] R. Tucker, M. Ruffini, L. Valcarenghi, D. R. Campelo, D. Simeonidou, L. Du, M. Marinescu, C. Middleton, S. Yin, T. Forde, K. Bourg, E. Dai, E. Harstead, P. Chanclou, H. Roberts, V. Jungnickel, S. Figuerola, T. Takahara, R. Yadav, P. Vetter, D. A. Khotimsky, and J. S. Wey, "Connected OFCity: Technology innovations for a smart city project," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 9, no. 2, pp. A245–A255, Feb 2017.

[118] *Network Functions Virtualisation (NFV); management and orchestration*, ETSI GS NFV-MAN GS NFV-MAN 001 V1.1.1, Jul. 2017.

[119] E. Riccardi, P. Gunning, G. de Dios, M. Quagliotti, V. Lopez, and A. Lord, "An Operator view on the Introduction of White Boxes into Optical Networks," *Journal of Lightwave Technology*, vol. 36, no. 15, pp. 3062–3072, Aug 2018.

[120] G. Agrawal, *Fiber-Optic Communication Systems*. Wiley-Interscience, 1997.

[121] G. Bosco, V. Curri, A. Carena, P. Poggiolini, and F. Forghieri, "On the Performance of Nyquist-WDM Terabit Superchannels Based on PM-BPSK, PM-QPSK, PM-8QAM or PM-16QAM Subcarriers," *Journal of Lightwave Technology*, vol. 29, no. 1, pp. 53–61, Jan 2011.

[122] N. Sambo, M. Secondini, F. Cugini, G. Bottari, P. Iovanna, F. Cavaliere, and P. Castoldi, "Modeling and Distributed Provisioning in 10-40-100-Gb/s Multirate Wavelength Switched Optical Networks," *Journal of Lightwave Technology*, vol. 29, no. 9, pp. 1248–1257, May 2011.

[123] L. Valcarenghi, F. Giannone, D. Manicone, and P. Castoldi, "Virtualized eNB latency limits," in *2017 19th International Conference on Transparent Optical Networks (ICTON)*, July 2017, pp. 1–4.

[124] K. Christodoulopoulos, N. Sambo, N. Argyris, P. Giardina, G. Kanakis, A. Kretsis, F. Fresi, A. Sgambelluri, G. Bernini, C. Delezoide, F. Cugini,

H. Avramopoulos, and E. Varvarigos, "Observe-Decide-Act: Experimental Demonstration of a Self-Healing Network," in *2018 Optical Fiber Communications Conference and Exposition (OFC)*, March 2018, pp. 1–3.

[125] M. Dallaglio, N. Sambo, F. Cugini, and P. Castoldi, "YANG Models for Vendor-Neutral Optical Networks, Reconfigurable through State Machine," *IEEE Communications Magazine*, vol. 55, no. 8, pp. 170–178, Aug 2017.

[126] N. Sambo, A. Giorgetti, F. Cugini, and P. Castoldi, "Sliceable Transponders: Pre-Programmed OAM, Control, and Management," *Journal of Lightwave Technology*, vol. 36, no. 7, pp. 1403–1410, April 2018.

[127] H. Gupta, D. Manicone, F. Giannone, K. Kondepu, A. Franklin, P. Castoldi, and L. Valcarenghi, "How much is fronthaul latency budget impacted by RAN virtualisation ?" in *2017 IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN)*, Nov 2017, pp. 315–320.

[128] *Network Functions Virtualisation (NFV); Architectural Framework*, ETSI Group Specification GS NFV 002, Oct. 2013.

[129] A. de la Oliva, X. Li, X. Costa-Perez, C. J. Bernardos, P. Bertin, P. Iovanna, T. Deiss, J. Mangues, A. Mourad, C. Casetti, J. E. Gonzalez, and A. Azcorra, "5G-TRANSFORMER: Slicing and Orchestrating Transport Networks for Industry Verticals," *IEEE Communications Magazine*, vol. 56, no. 8, pp. 78–84, August 2018.

[130] *Network Functions Virtualisation (NFV) Management and Orchestration*, ETSI Group Specification GS NFV-MAN 001, Dec. 2014.

[131] *Network Functions Virtualisation (NFV) Release 3; Evolution and Ecosystem; Report on Network Slicing Support with ETSI NFV Architecture Framework*, ETSI Group Specification GR NFV-EVE 012, Dec. 2017.

[132] *Network Functions Virtualisation (NFV) Release 3; Management and Orchestration; Report on architecture options to support multiple administrative domains*, ETSI Group Specification GR NFV-IFA 028, Jan. 2018.

[133] Open Network Foundation. ONF. [Online]. Available: https://onosproject.org/

[134] M. Mahalingam, D. G. Dutt, K. Duda, P. Agarwal, L. Kreeger, T. Sridhar, M. Bursell, and C. Wright, "Virtual eXtensible Local Area Network (VXLAN): A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks," *RFC*, vol. 7348, pp. 1–22, 2014.